



# Spatial approximation of the radiation transport equation using a subgrid-scale finite element method

Matias Avila, Ramon Codina\*, Javier Principe

Universitat Politècnica de Catalunya, Jordi Girona 1-3, Edifici C1, 08034 Barcelona, Spain

## ARTICLE INFO

### Article history:

Received 8 August 2009

Received in revised form 22 October 2010

Accepted 2 November 2010

Available online 9 November 2010

### Keywords:

Radiative transfer equation

Stabilized finite element methods

SUPG

Orthogonal subscales

Discrete ordinates method

## ABSTRACT

In this paper we present stabilized finite element methods to discretize in space the monochromatic radiation transport equation. These methods are based on the decomposition of the unknowns into resolvable and subgrid scales, with an approximation for the latter that yields a problem to be solved for the former. This approach allows us to design the algorithmic parameters on which the method depends, which we do here when the discrete ordinates method is used for the directional approximation. We concentrate on two stabilized methods, namely, the classical SUPG technique and the orthogonal subscale stabilization. A numerical analysis of the *spatial* approximation for both formulations is performed, which shows that they have a similar behavior: they are both stable and optimally convergent in the same mesh-dependent norm. A comparison with the behavior of the Galerkin method, for which a non-standard numerical analysis is done, is also presented.

© 2010 Elsevier B.V. All rights reserved.

## 1. Introduction

Radiation is energy propagation due to movement of subatomic particles or photons. Thermal radiation in particular refers to radiation caused by electromagnetic waves or photons. From the mathematical point of view, the problem consists in finding the radiative intensity field  $u$ , which depends on the position ( $\mathbf{x}$ ), on the propagating direction ( $\mathbf{s}$ ) and on the frequency or wave length ( $\lambda$ ), that is  $u = u(\lambda, \mathbf{x}, \mathbf{s})$ . In many cases of interest it is possible to assume that different frequencies do not interact and therefore the radiative intensity is computed for each frequency separately. The problem is still very hard to approximate numerically, as it involves in one way or another the discretization of the spatial and directional domains to obtain approximations to the solution of the integro-differential radiative transport equation (RTE). An exception is the Monte Carlo method, in which these discretizations are not explicitly built but the movement of individual photons is followed invoking statistical concepts to solve the problem. It is widely recognized as a simple and efficient method, but also very time consuming, especially in three dimensions.

During the last decades there were numerous efforts to develop suitable numerical schemes for the radiative transfer equation [16]. Several options for the directional discretization have been considered in the literature, including the discrete ordinates method (DOM) and the method of spherical harmonics (also called PN

approximation). In both cases the directional discretization transforms the integro-differential RTE into a set of coupled differential equations. In the case of the DOM the unknown of each equation is the radiative intensity in a given direction and integration over the solid angle is replaced by a quadrature sum, which is a set of discrete ordinate directions and the corresponding weights. The DOM was first proposed by Chandrasekhar [4] in his work on stellar and atmospheric radiation, analyzing radiation problems within a plane parallel medium. The PN approximation, which will not be used here, was first proposed by Jeans [12] in his work on radiative transfer in stars.

Besides the selection of the directional discretization, a proper spatial discretization is needed. When the problem also involves convective heat transfer, the numerical scheme for the RTE should be compatible with those for the mass, momentum and energy conservation equations governing the flow field. Implementation of inhomogeneity and anisotropic radiative scattering is required for practical applications in a multidimensional complex geometry, using different grids. In these situations natural candidates are finite volume methods (FVM) and finite element methods (FEM). These methods have the advantage that can deal with complex geometric shape and boundary conditions conveniently without increasing the computational complexity.

A finite element approximation of radiative heat transfer in one-dimensional problems was proposed by Viskanta [18] in 1965. Razzaque et al. [15] studied the finite element solution of radiative heat transfer in a two dimensional rectangular enclosure. Fiveland [20] developed a finite element formulation based on the DOM to solve absorbing, emitting and isotropic scattering in

\* Corresponding author.

E-mail addresses: [mavila@cimne.upc.edu](mailto:mavila@cimne.upc.edu) (M. Avila), [ramon.codina@upc.edu](mailto:ramon.codina@upc.edu) (R. Codina), [principe@cimne.upc.edu](mailto:principe@cimne.upc.edu) (J. Principe).

multidimensional problems. That formulation, however, cannot deal with the problem of anisotropic scattering. Richling et al. [19] formulated the radiative transfer equation in three dimensions for discrete ordinates using finite elements for anisotropically scattering media.

It is well known that a careful numerical formulation of the problem is needed to avoid numerical instabilities due to the first order hyperbolic nature of the problem. The Galerkin formulation is known to be unstable and therefore a stabilized formulation is necessary. For example, Kanschat [13] applied and analyzed the streamline upwind Petrov Galerkin (SUPG) formulation [3] for the problem continuous in the propagating direction. A discontinuous Galerkin approximation in space combined with the DOM for the directional discretization was proposed in [1].

In this work we propose a stabilized finite element formulation based on an arbitrary angular discretization for solving the radiative transfer equation in multidimensional geometries for absorbing, emitting and anisotropic scattering media. We analyze the SUPG method and the orthogonal subscales stabilization (OSS) method [5], which can be described in the variational multiscale framework introduced in [11]. A comparison with the behavior of the Galerkin method, for which a non-standard numerical analysis is performed, is also presented.

The paper is organized as follows. After describing the problem in Section 2, we present in Section 3.1 the spatial discretization, which is based on the variational multiscale formulation and the algebraic approximation to the subscales, leaving the discussion of the choice of the stabilization parameters to Section 3.3.2. A general directional discretization and its particular form for the DOM is presented in Section 3.2. A complete numerical analysis of the formulation is presented in Section 4, where stability and optimal convergence of the SUPG and OSS methods are proved, together with a non-standard stability and convergence analysis of the Galerkin method. The accuracy and efficiency of the scheme are discussed in Section 5, where some numerical experiments are presented. Concluding remarks close the paper in Section 6.

Register for free at <https://www.scipedia.com> to download the version without the watermark

## 2. Problem statement

### 2.1. Boundary value problem

Let  $\Omega \subset \mathbb{R}^3$  and let  $S^2$  be the unit sphere in  $\mathbb{R}^3$ . For conciseness, we consider the three dimensional problem, but all what follows can be applied to the two-dimensional case as well.

The monochromatic radiative transfer problem consists in finding  $u : \Omega \times S^2 \rightarrow \mathbb{R}$  such that

$$Lu = f \quad \text{in } \Omega \times S^2, \quad (1)$$

where the source of intensity  $f(\mathbf{x}, \mathbf{s})$  is a given function (depending on the temperature in thermal radiation) and the operator  $L$  is defined as

$$Lu(\mathbf{x}, \mathbf{s}) = \mathbf{s} \cdot \nabla u(\mathbf{x}, \mathbf{s}) + \kappa(\mathbf{x})u(\mathbf{x}, \mathbf{s}) + S_\sigma u(\mathbf{x}, \mathbf{s}), \quad (\mathbf{x}, \mathbf{s}) \in \Omega \times S^2, \quad (2)$$

where

$$S_\sigma u(\mathbf{x}, \mathbf{s}) := \sigma_s(\mathbf{x})S_1 u(\mathbf{x}, \mathbf{s}), \quad (3)$$

$$S_1 u(\mathbf{x}, \mathbf{s}) := u(\mathbf{x}, \mathbf{s}) - \frac{1}{4\pi} \int_{S^2} \phi(\mathbf{s}, \mathbf{s}') u(\mathbf{x}, \mathbf{s}') d\mathbf{s}'. \quad (4)$$

For clarity, the arguments on which the functions depend have been explicitly displayed. Functions  $\kappa(\mathbf{x}) \geq 0$  and  $\sigma_s(\mathbf{x}) \geq 0$  in (2) and (3) are the absorption and extinction coefficients, respectively. They only need to be bounded for the following developments, although we will consider them constant in the numerical analysis for simplicity.

The operator  $S_\sigma$  defined in (3) is the so called scattering operator. It depends on the phase function  $\phi \in C^\infty(S^2 \times S^2; \mathbb{R}^+)$ , which is normalized in such a way that

$$\int_{S^2} \phi(\mathbf{s}, \mathbf{s}') d\mathbf{s}' = 4\pi \quad \forall \mathbf{s} \in S^2.$$

According to the physical model, the phase function  $\phi(\mathbf{s}, \mathbf{s}')$  usually depends only on the cosine of the angle between  $\mathbf{s}$  and  $\mathbf{s}'$ . For heterogeneous media, it could also depend on the position  $\mathbf{x}$ , although we will not consider this possibility here.

The boundary  $\Gamma = \partial\Omega \times S^2$  of  $\Omega \times S^2$  is divided into the inflow  $\Gamma^-$  and outflow  $\Gamma^+$  boundaries, defined as

$$\Gamma^- = \{(\mathbf{x}, \mathbf{s}) \in \Gamma | \mathbf{s} \cdot \mathbf{n} < 0\}, \quad \Gamma^+ = \{(\mathbf{x}, \mathbf{s}) \in \Gamma | \mathbf{s} \cdot \mathbf{n} \geq 0\}, \quad (5)$$

where  $\mathbf{n}$  is the unit normal vector pointing outwards  $\partial\Omega$  at  $\mathbf{x}$ . We shall also make use of the hemispheres

$$S^- := \{\mathbf{s} \in S^2 | \mathbf{s} \cdot \mathbf{n} < 0\}, \quad S^+ := \{\mathbf{s} \in S^2 | \mathbf{s} \cdot \mathbf{n} \geq 0\}, \quad (6)$$

which are defined for each  $\mathbf{x} \in \partial\Omega$ .

For simplicity, in the description and analysis of the formulation we shall supply (1) with the simplest boundary condition  $u = 0$  on  $\Gamma^-$ , although in the numerical examples we shall deal also with emissive and reflective boundary conditions. Changes required to extend the numerical approximation and its analysis to this situation are explained in Section 2.3.

### 2.2. Variational form

In order to write the weak form of the problem let us introduce the spaces

$$\mathcal{V} = \{u : \Omega \times S^2 \rightarrow \mathbb{R} | u, \mathbf{s} \cdot \nabla u \in L^2(\Omega) \forall \mathbf{s} \in S^2\},$$

$$\mathcal{W} = \{u : \Omega \times S^2 \rightarrow \mathbb{R} | \|u\|_\Omega, \|\mathbf{s} \cdot \nabla u\|_\Omega \in L^2(S^2)\} = L^2(S^2; \mathcal{V}),$$

where  $\|u\|_\Omega$  is the usual  $L^2(\Omega)$ -norm. We also define the inner product for functions in  $\mathcal{W}$  (not the one associated to its topology) as

$$(u, v) = \int_{S^2} \int_\Omega u(\mathbf{x}, \mathbf{s}) v(\mathbf{x}, \mathbf{s}) d\mathbf{x} d\mathbf{s} = \int_{S^2} (u, v)_\Omega d\mathbf{s}, \quad (7)$$

where  $(u, v)_\Omega$  is the usual  $L^2(\Omega)$ -inner product. The norm associated to  $(\cdot, \cdot)$  is written as  $\|u\| = (u, u)^{1/2}$ . If  $\gamma \subset \partial\Omega \times S^2$ , we define

$$(u, v)_\gamma = \int_\gamma u(\mathbf{x}, \mathbf{s}) v(\mathbf{x}, \mathbf{s}) |\mathbf{n} \cdot \mathbf{s}| d\mathbf{x} d\mathbf{s}$$

and the associated norm  $\|u\|_\gamma = (u, u)_\gamma^{1/2}$ . In particular, we will use this definition for  $\gamma = \Gamma^+$ , case in which  $|\mathbf{n} \cdot \mathbf{s}| = \mathbf{n} \cdot \mathbf{s}$ .

The weak form of problem (1) consists in finding  $u \in \mathcal{W}$  such that

$$\begin{aligned} \mathcal{B}(u, v) &:= (Lu, v) = (\mathbf{s} \cdot \nabla u, v) + (\kappa u, v) + (S_\sigma u, v) = (f, v) \\ &=: \mathcal{L}(v) \quad \forall v \in L^2(S^2; L^2(\Omega)). \end{aligned} \quad (8)$$

By assumption,  $\phi$  is a bounded and symmetric function on  $S^2 \times S^2$ , and therefore  $S_1$  defined in (4) is a self-adjoint operator from  $L^2(S^2)$  onto itself. It is a compact perturbation of the identity operator and has a real and countable spectrum confined to the interval  $[0, 1]$ . The set of eigenfunctions corresponding to the eigenvalue  $\lambda_0 = 0$  contains at least the constants on  $S^2$  and, furthermore, zero is an isolated eigenvalue.

By the Hilbert–Schmidt theorem there is an orthonormal set  $\{u_n(\mathbf{s})\}$  of eigenfunctions corresponding to non-zero eigenvalues  $\{\lambda_n\}$  of  $S_1$  such that every function  $u(\mathbf{x}, \cdot) \in L^2(S^2)$  has a unique decomposition of the form

$$u(\mathbf{x}, \mathbf{s}) = \sum_{n=1}^{\infty} \alpha_n(\mathbf{x}) u_n(\mathbf{s}) + u_0(\mathbf{x}, \mathbf{s}), \quad (9)$$

where  $\alpha_n(\mathbf{x}) \in \mathbb{R}$ ,  $u_0(\mathbf{x}, \mathbf{s}) \in \ker S_1$  and

$$S_1 u(\mathbf{x}, \mathbf{s}) = \sum_{n=1}^{\infty} \lambda_n \alpha_n(\mathbf{x}) u_n(\mathbf{s}). \quad (10)$$

For  $\mathbf{x} \in \Omega$  fixed,  $L^2(\mathcal{S}^2)$  can be decomposed as  $L^2(\mathcal{S}^2) = \ker S_1 \oplus (\ker S_1)^\perp$ . If we denote by  $\Pi^\perp$  the orthogonal projection onto  $(\ker S_1)^\perp$ , from (9) and (10) it is easily seen that the scattering operator satisfies the following properties when applied to functions  $u, v \in L^2(\Omega \times \mathcal{S}^2)$ :

$$\begin{aligned} (S_1 u, v) &= (\Pi^\perp S_1 u, \Pi^\perp v) = (\Pi^\perp S_1 v, \Pi^\perp u), \\ \lambda_1 \|\sigma_s \Pi^\perp u\| &\leq \|S_\sigma u\| \leq \|\sigma_s \Pi^\perp u\|, \\ \|S_1 u\|^2 &\leq (S_1 u, u). \end{aligned} \quad (11)$$

### 2.3. Emissive and reflective boundary conditions

A generalization of the homogeneous boundary conditions that we will use for the presentation of the method is a boundary condition of the form

$$\begin{aligned} u(\mathbf{x}, \mathbf{s})|_{\mathcal{S}_x^-} &= \epsilon \bar{u}(\mathbf{x}, \mathbf{s})|_{\mathcal{S}_x^-} + \frac{1-\epsilon}{\pi} A^+ u(\mathbf{x}), \\ A^\pm u(\mathbf{x}) &:= \int_{\mathcal{S}_x^\pm} u(\mathbf{x}, \mathbf{s}') |\mathbf{n} \cdot \mathbf{s}'| d\mathbf{s}', \end{aligned} \quad (12)$$

where  $\bar{u}$  is a given function and  $0 \leq \epsilon \leq 1$ . Note that the last term in (12) does not depend on  $\mathbf{s}$ . Such type of boundary condition corresponds to the so called emissive and reflective boundaries. In this case,  $\rho := 1 - \epsilon$  and  $\epsilon$  are the diffuse reflection and emissive wall coefficients, respectively, and  $\bar{u}(\mathbf{x}, \mathbf{s}) = I_b(\mathbf{x})$  is the blackbody radiation, given by  $I_b = \sigma_B T^4 / \pi$ , where  $T$  is the wall temperature and  $\sigma_B = 5.6704 \cdot 10^{-8} \text{ W/m}^2 \text{ K}$  is the Stefan–Boltzmann constant. Note that in this case the whole right-hand-side in (12) is independent of  $\mathbf{s}$ .

Condition (12) can be imposed weakly, leading to the following variational form of the problem: find  $u \in \mathcal{W}$  such that

$$\mathcal{B}(u, v) + \eta(u, v)_{\Gamma^-} - \eta \frac{1-\epsilon}{\pi} (A^+ u, v)_{\Gamma^-} = \mathcal{L}(v) + \eta(\epsilon \bar{u}, v)_{\Gamma^-} \quad \forall v \in \mathcal{W}, \quad (13)$$

where  $\eta$  is a dimensionless parameter that has to satisfy the conditions indicated later. Observe that the space of test functions is now  $\mathcal{W}$ , since traces of this test functions are required on  $\Gamma^-$ .

The introduction of the new terms in (13) with respect to (8) does not offer any problem for the numerical approximation. In fact, it can be understood that the boundary condition (12) is imposed with Nitsche's method [17].

The stability properties of both the continuous problem (13) and of its discrete finite element approximation are similar to those of (8). This is so because the positivity of the bilinear form in the left-hand-side of (13) can be guaranteed as explained in the following. Let us start noting that

$$\begin{aligned} A^\pm u(\mathbf{x}) &\leq \left( \int_{\mathcal{S}_x^\pm} |\mathbf{n} \cdot \mathbf{s}| d\mathbf{s} \right)^{1/2} \left( \int_{\mathcal{S}_x^\pm} u^2(\mathbf{x}, \mathbf{s}) |\mathbf{n} \cdot \mathbf{s}| d\mathbf{s} \right)^{1/2} \\ &= \sqrt{\pi} \left( \int_{\mathcal{S}_x^\pm} u^2(\mathbf{x}, \mathbf{s}) |\mathbf{n} \cdot \mathbf{s}| d\mathbf{s} \right)^{1/2} \end{aligned}$$

and therefore

$$A^+ u(\mathbf{x}) A^- u(\mathbf{x}) \leq \frac{\pi}{2} \left( \frac{1}{\alpha} \int_{\mathcal{S}_x^-} u^2(\mathbf{x}, \mathbf{s}) |\mathbf{n} \cdot \mathbf{s}| d\mathbf{s} + \alpha \int_{\mathcal{S}_x^+} u^2(\mathbf{x}, \mathbf{s}) |\mathbf{n} \cdot \mathbf{s}| d\mathbf{s} \right),$$

for all  $\alpha > 0$ . Using the positive definiteness of  $S_1$ , the fact that  $\kappa \geq 0$  and this last inequality it follows that

$$\begin{aligned} \mathcal{B}(u, u) + \eta(u, u)_{\Gamma^-} - \eta \frac{1-\epsilon}{\pi} (A^+ u, u)_{\Gamma^-} \\ \geq \frac{1}{2} (u, u)_{\Gamma^+} - \frac{1}{2} (u, u)_{\Gamma^-} + \eta(u, u)_{\Gamma^-} \\ - \eta \frac{1-\epsilon}{\pi} \int_{\partial\Omega} A^+ u(\mathbf{x}) A^- u(\mathbf{x}) d\mathbf{x} \\ \geq \|u\|_{\Gamma^+}^2 \left( \frac{1}{2} - \eta \frac{1-\epsilon}{2} \alpha \right) + \|u\|_{\Gamma^-}^2 \left( -\frac{1}{2} + \eta - \eta \frac{1-\epsilon}{2\alpha} \right). \end{aligned} \quad (14)$$

Positivity of (14) follows from conditions

$$\eta \leq \frac{1}{\alpha(1-\epsilon)}, \quad \eta \left( 1 - \frac{1-\epsilon}{2\alpha} \right) \geq \frac{1}{2}. \quad (15)$$

For  $\epsilon = 1$  the only requirement is  $\eta \geq \frac{1}{2}$ , whereas for  $\epsilon < 1$  it is easily checked that the above conditions are feasible. For example, for  $\alpha = 1 - \epsilon$  they read  $1 \leq \eta \leq (1 - \epsilon)^{-2}$ . In fact, a little analysis shows that the upper bound for  $\eta$  is maximized if  $\alpha = \frac{1}{\rho} - \sqrt{\frac{1}{\rho^2} - 1}$  (where  $\rho = 1 - \epsilon$ ), which satisfies  $\alpha \geq \frac{\rho}{2}$  and is obtained by choosing  $\alpha$  as the smallest value satisfying  $(\rho\alpha)^{-1} \geq [2(1 - \frac{\rho}{2\alpha})]^{-1}$ . Note that when the  $=$  sign holds in the first inequality of (15), control on  $\|u\|_{\Gamma^+}$  is lost. If this norm is to be included in the stability norm (see (42) below), the  $\leq$  symbol has to be replaced by  $<$ . In the case  $\epsilon = 0$  one must choose  $\alpha = 1$  and thus  $\eta = 1$ . Therefore,  $\epsilon > 0$  is needed if (15) has to hold with strict inequalities.

The positivity of (14) is enough to extend the formulation and analysis of the methods that follow to (12) with minor modifications.

## 3. Numerical approximation

In this section we consider the numerical approximation of problem (8). Spatial and directional discretizations will be considered independently. Our main concern is the former, and we will therefore consider a generic finite dimensional space of functions defined on  $\mathcal{S}^2$ . However, we will also particularize our formulation with regard to the latter, with the aim of performing numerical experiments of Section 5.

### 3.1. Spatial discretization

Let us consider a finite element partition  $\mathcal{P}_h = \{K\}$  of the domain  $\Omega$  of diameter  $h$ . From this finite element partition we build up conforming finite element spaces  $\mathcal{V}_h \subset \mathcal{V}$  in the usual manner. Let also  $\mathcal{W}_h = L^2(\mathcal{S}^2; \mathcal{V}_h)$ . Discrete test functions will also be taken in this space.

#### 3.1.1. Galerkin finite element approximation

The spatial Galerkin finite element approximation of problem (8) consists in finding  $u_h \in \mathcal{W}_h$  such that

$$\mathcal{B}(u_h, v_h) = \mathcal{L}(v_h) \quad \forall v_h \in \mathcal{W}_h. \quad (16)$$

The question that arises once the discrete problem is set is whether it is stable or not. The bilinear form  $\mathcal{B}(u_h, v_h)$  is not coercive with respect to the graph norm (in particular, with respect to the derivatives involved), as we shall see later. In the particular case of null scattering ( $\sigma_s = 0$  on  $\Omega$ ), the problem decouples into a system of convection–reaction equations on  $\Omega$ . This system is hyperbolic, and the Galerkin finite element method is known to produce spurious oscillations in this case.

#### 3.1.2. Stabilized finite element approximation using subscales

In this section we describe the finite element approximation proposed, which can be cast in the variational multiscale framework proposed in [11]. For completeness, we briefly describe it in the following, also adapted to our particular proposal.

As in [11], let us split the continuous space  $\mathcal{V}$  as  $\mathcal{V} = \mathcal{V}_h \oplus \tilde{\mathcal{V}}$ , where  $\tilde{\mathcal{V}}$  is any space to complete  $\mathcal{V}_h$  in  $\mathcal{V}$ , obviously infinite-dimensional. From  $\mathcal{V}_h$  and  $\tilde{\mathcal{V}}$  we may define  $\mathcal{W}_h = L^2(\mathcal{S}^2; \mathcal{V}_h)$  and  $\tilde{\mathcal{W}} = L^2(\mathcal{S}^2; \tilde{\mathcal{V}})$ , which satisfy  $\mathcal{W} = \mathcal{W}_h \oplus \tilde{\mathcal{W}}$ . Since  $\tilde{u} \in \tilde{\mathcal{W}}$  represents the component of  $u$  whose spatial dependence cannot be represented in the finite element space, we call  $\tilde{\mathcal{W}}$  the space of *subcales* or *subgrid scales*. We may also decompose the space of test functions as  $L^2(\mathcal{S}^2; L^2(\Omega)) = \mathcal{W}_h + \tilde{\mathcal{W}}_T$ .

The weak form of the continuous problem (8) is exactly equivalent to find  $u_h \in \mathcal{W}_h$  and  $\tilde{u} \in \tilde{\mathcal{W}}$  such that

$$\mathcal{B}(u_h, v_h) + \mathcal{B}(\tilde{u}, v_h) = \mathcal{L}(v_h) \quad \forall v_h \in \mathcal{W}_h, \quad (17)$$

$$\mathcal{B}(u_h, \tilde{v}) + \mathcal{B}(\tilde{u}, \tilde{v}) = \mathcal{L}(\tilde{v}) \quad \forall \tilde{v} \in \tilde{\mathcal{W}}_T. \quad (18)$$

Integrating the convective term by parts in the second term in the left-hand-side of (17) and in (18) it is found that problem (17) and (18) can be written as

$$\mathcal{B}(u_h, v_h) + (L^* v_h, \tilde{u}) + (v_h, \tilde{u})_{\Gamma^+} = (v_h, f), \quad (19)$$

$$(\tilde{v}, L\tilde{u}) = (\tilde{v}, f - Lu_h), \quad (20)$$

where  $L^*$  denotes the adjoint operator, which is defined as

$$L^* u(\mathbf{x}, \mathbf{s}) = -\mathbf{s} \cdot \nabla u(\mathbf{x}, \mathbf{s}) + \kappa u(\mathbf{x}, \mathbf{s}) + S_\sigma u(\mathbf{x}, \mathbf{s}).$$

Eq. (20) is equivalent to

$$L\tilde{u} = f - Lu_h + v_{h,\text{ort}} \quad \text{in } \Omega, \quad v_{h,\text{ort}} \in \tilde{\mathcal{W}}_T^\perp, \quad (21)$$

where  $v_{h,\text{ort}}$  is responsible to enforce that the previous equation holds in the space of the subscales. The goal of all subscale methods is to approximate  $\tilde{u}$  to end up with a modified problem for  $u_h$  with enhanced stability properties.

There are several possibilities to deal with problem (21). We consider the algebraic approximation  $L\tilde{u} \approx \tau^{-1}\tilde{u}$ , which has to be understood in the  $L^2$ -norm (see [6]). Another possibility is to approximate  $\tilde{u}$  by bubble functions, which leads to a similar method for a proper identification of this bubble functions (see [2], for example). When replaced into Eq. (21) this approximation gives

$$\tilde{u} = \tau(f - Lu_h + v_{h,\text{ort}}), \quad (22)$$

where  $\tau$  is an algorithmic parameter depending on the geometry of each element domain  $K$  and the coefficients of operator  $L$ . This approximation for  $\tilde{u}$  is intended to mimic the effect of the exact subscales in the volume integral of (19), whereas the integral over the boundary  $\Gamma^+$  will be neglected. It remains to define the stabilization parameter  $\tau$  in terms of the equation coefficients and the mesh size and to define  $v_{h,\text{ort}}$ , thus selecting the space of subscales. The approximation performed to obtain  $\tau$  is based on an (approximate) Fourier analysis of the problem as in [6] and will be discussed in Section 3.3.2 after the directional discretization is introduced. The choice of the space of subscales is discussed in the rest of this section, where two different possibilities are considered.

**3.1.2.1. Algebraic subscales and the SUPG method.** The simplest choice for the space of subscales is to take  $v_{h,\text{ort}} = 0$ , which implies that the subscales belong to the space of the residuals. This results in what we will call algebraic subgrid scale method (ASGS). Inserting (22) in (19), the discrete problem in this case reads: find  $u_h \in \mathcal{W}_h$  such that

$$\mathcal{B}_{\text{asgs}}(u_h, v_h) = \mathcal{L}_{\text{asgs}}(v_h) \quad \forall v_h \in \mathcal{W}_h,$$

where

$$\mathcal{B}_{\text{asgs}}(u_h, v_h) := \mathcal{B}(u_h, v_h) + (-L^* v_h, \tau Lu_h),$$

$$\mathcal{L}_{\text{asgs}}(v_h) := \mathcal{L}(v_h) + (-L^* v_h, \tau f).$$

This method is a generalization of the well known SUPG method in which only the advective term is considered to weight the subscales

instead of the whole adjoint operator  $L^*$ . In this method the discrete problem reads: find  $u_h \in \mathcal{W}_h$  such that

$$\mathcal{B}_{\text{supg}}(u_h, v_h) = \mathcal{L}_{\text{supg}}(v_h) \quad \forall v_h \in \mathcal{W}_h, \quad (23)$$

where

$$\mathcal{B}_{\text{supg}}(u_h, v_h) := \mathcal{B}(u_h, v_h) + (\mathbf{s} \cdot \nabla v_h, \tau Lu_h), \quad (24)$$

$$\mathcal{L}_{\text{supg}}(v_h) := \mathcal{L}(v_h) + (\mathbf{s} \cdot \nabla v_h, \tau f).$$

**3.1.2.2. Orthogonal subscales.** The starting point has been the decomposition  $\mathcal{V} = \mathcal{V}_h \oplus \tilde{\mathcal{V}}$ . Among the possibilities to choose  $\tilde{\mathcal{V}}$ , a particular choice is to take the space for the subscales orthogonal to the finite element space, that is to say,

$$\tilde{\mathcal{V}} = \mathcal{V}_h^\perp \cap \mathcal{V} \approx \mathcal{V}_h^\perp, \quad (25)$$

where the symbol  $\approx$  has to be understood in the sense that conformity is violated. Using the same reasoning as in [6], this approximation together with some additional simplifications lead to the expression for the subscales

$$\tilde{u} = -\tau P_h^\perp (\mathbf{s} \cdot \nabla u_h),$$

with  $P_h^\perp = I - P_h$ ,  $I$  being the identity and  $P_h$  the  $L^2$ -projection onto the finite element space. Replacing this expression into Eq. (19) and taking into account that subscale functions vanish on  $\partial\Omega$  we get the final discrete problem: find  $u_h \in \mathcal{W}_h$  such that

$$\mathcal{B}_{\text{oss}}(u_h, v_h) = \mathcal{L}(v_h) \quad \forall v_h \in \mathcal{W}_h, \quad (26)$$

where, for constant physical coefficients and uniform meshes,

$$\mathcal{B}_{\text{oss}}(u_h, v_h) := \mathcal{B}(u_h, v_h) + (\mathbf{s} \cdot \nabla v_h, \tau P_h^\perp (\mathbf{s} \cdot \nabla u_h)). \quad (27)$$

The simplifying assumptions have yielded a method that is easy to implement and with good stability properties, as we shall see in Section 4. Comparing expressions (23) and (27) it is observed that the latter has less terms to evaluate, since it is not necessary to weight the whole residual by  $\mathbf{s} \cdot \nabla v_h$ . However, the term  $P_h^\perp (\mathbf{s} \cdot \nabla u_h)$  leads to a wider connectivity between the mesh nodes. Nevertheless, iterative schemes may be devised to deal with this coupling, which may be very effective when the RTE is coupled with nonlinear flow problems.

## 3.2. Directional discretization

### 3.2.1. Approximation of the directional component

As mentioned in Section 1, there exist several possible choices for the directional discretization. Introducing a generic basis  $\{\psi^\alpha(\mathbf{s}), \alpha = 1, \dots, N\}$  to approximate  $L^2(\mathcal{S}^2)$  with a space of dimension  $N$ , we can approximate  $\mathcal{W} = L^2(\mathcal{S}^2; \mathcal{V})$  by

$$\mathcal{W}_N := \left\{ v \in L^2(\mathcal{S}^2; \mathcal{V}) \mid v(\mathbf{x}, \mathbf{s}) = \sum_{\alpha=1}^N \psi^\alpha(\mathbf{s}) v^\alpha(\mathbf{x}), v^\alpha(\mathbf{x}) \in \mathcal{V} \quad \forall \alpha \right\}.$$

The space of test functions  $L^2(\mathcal{S}^2; L^2(\Omega))$  can be approximated similarly, replacing the condition  $v^\alpha(\mathbf{x}) \in \mathcal{V}$  by  $v^\alpha(\mathbf{x}) \in L^2(\Omega)$ . The resulting space is denoted by  $\mathcal{W}_{T,N}$ .

The Galerkin method applied to the directional discretization of problem (8) consists in finding  $u_N \in \mathcal{W}_N$  such that

$$\mathcal{B}(u_N, v_N) = \mathcal{L}(v_N) \quad \forall v_N \in \mathcal{W}_{T,N}. \quad (28)$$

If  $u^\alpha$  and  $v^\alpha$  are the components of the unknown and test function in  $\mathcal{W}_N$ , we may write

$$\begin{aligned} \mathcal{B}(u_N, v_N) &= \sum_{\alpha, \beta=1}^N \left[ \left( v^\alpha, A_i^{\alpha\beta} \partial_i u^\beta \right)_\Omega + \left( v^\alpha, S^{\alpha\beta} u^\beta \right)_\Omega \right], \\ \mathcal{L}(v_N) &= \sum_{\alpha=1}^N \left( v^\alpha, f^\alpha \right)_\Omega, \end{aligned} \quad (29)$$



where

$$A_i^{\alpha\beta} = \int_{S^2} s_i \psi^\alpha(\mathbf{s}) \psi^\beta(\mathbf{s}) d\mathbf{s}, \quad (30)$$

$$S^{\alpha\beta}(\mathbf{x}) = (\kappa(\mathbf{x}) + \sigma_s(\mathbf{x})) \int_{S^2} \psi^\alpha(\mathbf{s}) \psi^\beta(\mathbf{s}) d\mathbf{s} - \frac{\sigma_s(\mathbf{x})}{4\pi} \int_{S^2} \int_{S^2} \psi^\alpha(\mathbf{s}) \phi(\mathbf{s}, \mathbf{s}') \psi^\beta(\mathbf{s}') d\mathbf{s} d\mathbf{s}', \quad (31)$$

$$f^\alpha(\mathbf{x}) = \int_{S^2} \psi^\alpha(\mathbf{s}) f(\mathbf{x}, \mathbf{s}) d\mathbf{s}. \quad (32)$$

Repeated indexes  $i$  in (29) and below that run over the space dimensions (from 1 to 3) imply summation.

Defining the vector fields  $\mathbf{u} \in \mathcal{V}^N$  and  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^N$  by their components  $u^\alpha, f^\alpha, \alpha = 1, \dots, N$ , and introducing the vector operator  $\mathbf{L}$  defined by

$$[\mathbf{L}\mathbf{u}]^\alpha = \sum_{\beta=1}^N \left( A_i^{\alpha\beta} \partial_i u^\beta + S^{\alpha\beta} u^\beta \right), \quad (33)$$

we may write the discrete problem (28) (see (29)) as: find  $\mathbf{u} \in \mathcal{V}^N$  such that

$$(\mathbf{v}, \mathbf{L}\mathbf{u})_\Omega = (\mathbf{v}, \mathbf{f})_\Omega \quad \forall \mathbf{v} \in L^2(\Omega)^N. \quad (34)$$

### 3.2.2. The discrete ordinates method

The simplest way to discretize the directional domain  $S^2$  is using the so called *discrete ordinates method* (DOM). Let us consider a partition of the unit sphere of the form  $S^2 = \bigcup_{\alpha=1}^N S^{2,\alpha}$ , and let  $w_\alpha := \text{meas}(S^{2,\alpha})$ . The DOM is defined by setting  $\psi^\alpha(\mathbf{s}) = \chi^\alpha(\mathbf{s})$ , the characteristic function of  $S^{2,\alpha}$  (equal to 1 if  $\mathbf{s} \in S^{2,\alpha}$ , 0 otherwise).

If  $\mathbf{s}^\alpha$  is the barycenter of  $S^{2,\alpha}$ , expressions (30)–(32) may be approximated by

$$A_i^{\alpha\beta} = s_i^\alpha w_\alpha \delta_{\alpha\beta}, \quad (35)$$

$$S^{\alpha\beta}(\mathbf{x}) = (\kappa(\mathbf{x}) + \sigma_s(\mathbf{x})) w_\alpha \delta_{\alpha\beta} - \frac{\sigma_s(\mathbf{x})}{4\pi} w_\alpha w_\beta \phi(\mathbf{s}^\alpha, \mathbf{s}^\beta), \quad (36)$$

$$f^\alpha(\mathbf{x}) = w_\alpha f(\mathbf{x}, \mathbf{s}^\alpha). \quad (37)$$

In this case the unknown  $u^\alpha(\mathbf{x})$  represents the radiation intensity in the direction  $\mathbf{s}^\alpha$ , that is to say,  $u^\alpha(\mathbf{x}) = u(\mathbf{x}, \mathbf{s}^\alpha)$ .

### 3.3. Fully discrete problem using the discrete ordinates method

#### 3.3.1. Spatial and directional discretization

We may now proceed to the spatial discretization of problem (34). Let

$$\mathcal{W}_{h,N} := \left\{ v \in L^2(S^2; \mathcal{V}) \mid v(\mathbf{x}, \mathbf{s}) = \sum_{\alpha=1}^N \psi^\alpha(\mathbf{s}) v_h^\alpha(\mathbf{x}), v_h^\alpha(\mathbf{x}) \in \mathcal{V}_h \quad \forall \alpha \right\}$$

be the finite element space to approximate  $\mathcal{W}_N$ . This same space can be used to approximate the space of test functions  $\mathcal{W}_{T,N}$ . The Galerkin fully discrete problem, corresponding to the directional discretization of (16), consists in finding  $u_{h,N} \in \mathcal{W}_{h,N}$  such that

$$\mathcal{B}(u_{h,N}, v_{h,N}) = \mathcal{L}(v_{h,N}) \quad \forall v_{h,N} \in \mathcal{W}_{h,N}. \quad (38)$$

As explained earlier, this formulation lacks numerical stability. To design stabilized finite element methods we may proceed as in the previous subsection, simply replacing scalar-valued unknowns and test functions by their vector-valued counterparts, as well as the scalar operator  $L$  defined in (2) by the vector operator  $\mathbf{L}$  introduced in (33). In particular, the equation for the subscales  $\tilde{\mathbf{u}} \in \tilde{\mathcal{V}}^N$  will be

$$\tilde{P}(\tilde{\mathbf{L}}\tilde{\mathbf{u}}) = \tilde{P}(\mathbf{f} - \mathbf{L}\mathbf{u}_h), \quad (39)$$

where  $\tilde{P} = I$  for the ASGS formulation and  $\tilde{P} = P_h^\perp$  in the OSS method. It is understood that  $\tilde{P}$  acts componentwise. The approximation of  $\tilde{\mathbf{L}}\mathbf{u}$  we use is described next.

#### 3.3.2. The general approach to design the stabilization parameters

The unresolved subscales are modeled with the algebraic approximation in (22). The behavior of the stabilization parameter  $\tau$  can be analyzed using an approximate Fourier analysis of the problem, in the same way as it is done in [9,6].

Let us consider problem (39) posed in each element domain  $K$ . Our purpose is to approximate  $\tilde{\mathbf{L}}\mathbf{u} \approx \tau^{-1} \tilde{\mathbf{u}}$  in a certain sense, with  $\tau^{-1}$  a diagonal matrix that has to be determined and that we will call matrix of stabilization parameters. We propose to do this imposing that the induced  $L^2$ -norm of  $\tau^{-1}$  is an upper bound for the induced  $L^2$ -norm of  $\tilde{\mathbf{L}}$ , that is to say  $\|\tilde{\mathbf{L}}\|_{L^2(K)} \leq \|\tau^{-1}\|_{L^2(K)}$ . The symbol  $\leq$  has to be understood up to constants and holding independently of the equation coefficients. From an approximate Fourier analysis (see [9,6]) it may be concluded that  $\|\tilde{\mathbf{L}}\|_{L^2(K)} \leq |\hat{\tilde{\mathbf{L}}}(\mathbf{k}^0)|$  for a certain wave number, denoted  $\mathbf{k}^0$ , where the  $\hat{\tilde{\mathbf{L}}}$  is the algebraic operator resulting from the Fourier transform of  $\tilde{\mathbf{L}}\mathbf{u}$ . In view of this fact, our proposal is to choose  $\tau^{-1}$  such that  $|\hat{\tilde{\mathbf{L}}}(\mathbf{k}^0)| = |\tau^{-1}|$ . Obviously  $\mathbf{k}^0$  is unknown and has to be understood in this context as a vector of algorithmic coefficients.

The norm  $|\hat{\tilde{\mathbf{L}}}(\mathbf{k}^0)|$  can be computed as the square root of the maximum eigenvalue (in module) of the generalized eigenvalue problem

$$\hat{\tilde{\mathbf{L}}}(\mathbf{k}^0)^* \hat{\tilde{\mathbf{L}}}(\mathbf{k}^0) \mathbf{u} = \lambda \mathbf{u}. \quad (40)$$

This leads to an effective way to determine the expression of matrix of stabilization parameters  $\tau$ . Taking it as diagonal, it can be computed as  $\tau = \lambda_{\max}^{-1/2} \mathbf{I}$ .

The general idea exposed allows us to obtain the correct matrix of stabilization parameters for several problems (see [8,9] for an obtention of this matrix in the context of the hyperbolic wave equation and the three field formulation of the Stokes problem, for example). In particular, we will apply it now to the design of the  $\tau$  matrix for the DOM presented above. Let us also note that in some vector cases it is necessary to introduce a scaling matrix in the definition of  $\tau$ , which in the problem considered is not required.

#### 3.3.3. Tau matrix for the discrete ordinates method

For the sake of simplicity we will assume isotropic scattering, that is  $\phi = 1$ , and constant physical properties. Let us introduce the abbreviations  $\bar{\kappa} = \kappa + \sigma_s$ ,  $\bar{w}_\alpha = w_\alpha \frac{\sigma_s}{4\pi}$ . We wish to apply the previous ideas to problem (39). If we call  $\mathbf{r}$  the right-hand-side term, the Fourier transformed equation for the DOM reads

$$\sum_{\beta=1}^N w_\alpha \left[ i h^{-1}(\mathbf{k} \cdot \mathbf{s}^\alpha) \delta_{\alpha\beta} + \bar{\kappa} \delta_{\alpha\beta} - \bar{w}_\beta \right] \hat{u}^\beta = w_\alpha \hat{r}^\alpha,$$

where  $i = \sqrt{-1}$ ,  $\hat{r}^\alpha$  is the Fourier transform of  $\mathbf{r}(\mathbf{x}, \mathbf{s})$  evaluated at  $\mathbf{s}^\alpha$  and use has been made of (35)–(37). From this expression we see that we need to estimate the maximum eigenvalue of (40), the matrices involved having components

$$\hat{L}_{\alpha\beta}(\mathbf{k}^0) = i h^{-1}(\mathbf{k}^0 \cdot \mathbf{s}^\alpha) \delta_{\alpha\beta} + \bar{\kappa} \delta_{\alpha\beta} - \bar{w}_\beta,$$

$$\hat{L}_{\alpha\beta}^*(\mathbf{k}^0) = -i h^{-1}(\mathbf{k}^0 \cdot \mathbf{s}^\beta) \delta_{\alpha\beta} + \bar{\kappa} \delta_{\alpha\beta} - \bar{w}_\alpha.$$

After some algebraic manipulations that are omitted, it can be shown that

$$\lambda_{\max} \leq 2c \frac{\sigma_s}{N} + c^2 + \bar{\kappa}^2,$$

where  $c$  is a constant, independent of the equation coefficients and the number of modes in the DOM expansion  $N$ . In view of this, for

our problem we will take a diagonal matrix of stabilization parameters  $\tau$  given by

$$\tau = \left[ \left( \frac{c_1}{h} \right)^2 + (\kappa + \sigma_s)^2 + 2c_1 \frac{\sigma_s}{hN} \right]^{-1/2} \mathbf{I}, \quad (41)$$

where  $c_1$  is an algorithmic constant. We take it as  $c_1 = 2$  in the numerical examples of Section 5 using linear elements. In the analysis presented in the following section this constant is required to be large enough, but the value indicated is what we have found effective in practice.

The matrix of stabilization parameters is an algebraic approximation to the radiative transport operator in (1). The first term in (41) approximates the convective operator, the second term approximates the non integral reactive term and the last term approximates the integral operator. Observe that in the last term the directional and spatial discretization parameters  $N$  and  $h$  appear, and that it vanishes as  $N \rightarrow \infty$ . In any case, for practical values of  $N$  this last term is negligible, even if  $h$  decreases as  $N$  grows. A simple analysis of the magnitude of the different terms in (41) when both  $N \rightarrow \infty$  and  $h \rightarrow 0$  shows that the last term can never dominate. Thus, (41) has the behavior of the stabilization parameters usually found in the literature and that will be used in the numerical analysis presented next.

#### 4. Numerical analysis

##### 4.1. Preliminaries

In the present section we present a stability and convergence analysis of the SUPG and OSS methods, as well as a non-standard analysis of the Galerkin method, when they are used for the spatial discretization of the RTE, for simplicity with essential homogeneous boundary conditions on  $\Gamma^-$ .

For simplicity, we will consider quasi-uniform refinements, and thus all the element diameters can be bounded above and below by constants depending on  $h$ . The analysis of the SUPG and OSS methods for non quasi-uniform refinements and non uniform properties can be done using the strategy followed in [7]. We also consider uniform properties  $\kappa$  and  $\sigma_s$ . With these assumptions, the stabilization parameters can be considered constant in the whole computational domain.

For the SUPG and the OSS methods the norm in which the results will be presented is

$$\|v_h\|^2 := \tau \|\mathbf{s} \cdot \nabla v_h\|^2 + \|v_h\|_{\Gamma^+}^2 + \kappa \|v_h\|^2 + \lambda_1 \sigma_s \|\Pi^\perp v_h\|^2, \quad (42)$$

$$v_h(\cdot, \mathbf{s}) \in \mathcal{V}_h, \quad \mathbf{s} \in \mathcal{S}.$$

It is understood that the directional variable  $\mathbf{s}$  remains continuous, that is to say, the semi-discrete problem is analyzed, although minor modifications permit the extension of our analysis to the fully discrete DOM, for example. When the directional dependence is accounted for,  $v_h \in \mathcal{W}_h = L^2(\mathcal{S}^2; \mathcal{V}_h)$ . If  $\omega$  is a spatial domain, we will use the abbreviation  $\|v_h\|_\omega := \|v_h\|_{L^2(\mathcal{S}^2; L^2(\omega))}$ . Likewise, if  $|\cdot|_{H^i(\Omega)}$  is the seminorm of  $H^i(\Omega)$ , we will write  $\|v\|_{H^i(\Omega)} = |v|_{L^2(\mathcal{S}^2; H^i(\Omega))}$ .

Let  $u \in \mathcal{W}$  be the solution of the continuous problem and  $\hat{u}_h \in \mathcal{W}_h$  a finite element interpolant of degree  $p$ . If  $\|u - \hat{u}_h\|_{L^2(\mathcal{S}^2; H^i(\Omega))} \leq Ch^{p+1-i} |u|_{L^2(\mathcal{S}^2; H^{p+1}(\Omega))} =: \varepsilon_i(u)$ ,  $i = 0, 1$ , we will show that the error function of the SUPG and OSS methods in the norm (42) is given by

$$E(h) := \tau^{-1/2} \varepsilon_0(u) + \tau^{1/2} \varepsilon_1(u). \quad (43)$$

Obviously, we could express  $E(h)$  in terms of  $\varepsilon_0(u)$  or  $\varepsilon_1(u)$ , but we prefer to keep the explicit dependence on both to stress the behavior of  $E(h)$  in terms of the physical coefficients of the problem.

We will use the notation  $A_h \gtrsim B_h$  and  $A_h \lesssim B_h$  to indicate that  $A_h \geq CB_h$  and  $A_h \leq CB_h$ , respectively, where  $A_h$  and  $B_h$  are expressions that may depend on  $h$  and  $C$  is a generic constant, independent of  $h$  and of the physical parameters.

The expression of  $\tau$  we will use corresponds to the limit for  $N \rightarrow \infty$  of the diagonal in (41), that is to say,

$$\tau^{-2} = (\kappa + \sigma_s)^2 + c_1^2 h^{-2}. \quad (44)$$

Since the finite element partitions are assumed quasi-uniform, there is a positive constant  $C_{\text{inv}}$  independent of the mesh size  $h$  (the maximum of all element diameters), such that

$$\|\nabla v_h\|_K \leq C_{\text{inv}} h^{-1} \|v_h\|_K, \quad (45)$$

for all finite element functions  $v_h$  defined on an element  $K \in \mathcal{P}_h$ . Remember that the subscript  $K$  denotes that the spatial integral involved in  $\|\cdot\|$  is carried over element  $K$ . Similarly, the trace inequality

$$\|v\|_{\partial K}^2 \leq C_{\text{trace}} \left( h^{-1} \|v\|_K^2 + h \|\nabla v\|_K^2 \right) \quad (46)$$

is assumed to hold for functions  $v \in L^2(\mathcal{S}^2; H^1(K))$ ,  $K \in \mathcal{P}_h$ . Now, subscript  $\partial K$  denotes the  $L^2(\partial K)$ -norm. The last term can be dropped if  $v$  is a polynomial on the element domain  $K$ . Thus, if  $v_h$  is a piecewise continuous polynomial, it follows that

$$\|v_h\|_{\partial K}^2 \leq C_{\text{trace}} h^{-1} \|v_h\|_K^2. \quad (47)$$

**Lemma 1.** For sufficiently smooth solutions  $u$  of the continuous problem there holds

$$\|u - \hat{u}_h\|_{\Gamma^+} \leq E(h).$$

**Proof.** The results follows easily from the definition of  $\|\cdot\|$  and the previous assumptions, together with the trace inequality (46), which implies

$$\|u - \hat{u}_h\|_{\Gamma^+} \lesssim h^{-1/2} \varepsilon_0(u) + h^{1/2} \varepsilon_1(u). \quad (48)$$

This inequality will be used later on.  $\square$

**Remark 1.** The previous result makes sense for smooth enough functions  $u$ , in particular for  $u$  at least in  $L^2(\mathcal{S}^2; H^1(\Omega))$ . Thus, the traces of  $u$  on  $\Gamma^+$  are well defined for almost every direction, as well as the traces of functions in the finite element spaces we have constructed. However, if instead of seeking the order of convergence we only want to prove convergence towards a solution with the minimum regularity requirements, that is to say,  $u \in \mathcal{W}$ , the trace of  $u$  on  $\Gamma^+$  is not necessarily defined. As it is shown in [10] (Lemma 3.1) this trace makes sense if  $\Gamma^+$  and  $\Gamma^-$  defined in (5) are well separated.

##### 4.2. SUPG method

In this subsection we will prove that the solution of the discrete problem (23) is stable and convergent to the solution of the continuous problem (1).

**Lemma 2 (Coercivity of the SUPG method).** The bilinear form  $\mathcal{B}_{\text{supg}}$  defined in (24) satisfies

$$\mathcal{B}_{\text{supg}}(v_h, v_h) \gtrsim \|v_h\|^2 \quad \forall v_h \in \mathcal{W}_h.$$

**Proof.** From the definition of  $\mathcal{B}_{\text{supg}}$  it follows that

$$\begin{aligned}\mathcal{B}_{\text{supg}}(v_h, v_h) &= (\mathbf{s} \cdot \nabla v_h, v_h) + \kappa(v_h, v_h) + (S_\sigma v_h, v_h) \\ &\quad + \tau(\mathbf{s} \cdot \nabla v_h, \mathbf{s} \cdot \nabla v_h) + \kappa\tau(v_h, \mathbf{s} \cdot \nabla v_h) \\ &\quad + \tau(S_\sigma v_h, \mathbf{s} \cdot \nabla v_h) \geq \frac{1}{2}\|v_h\|_{r^+}^2 + \kappa\|v_h\|^2 \\ &\quad + \sigma_s(S_1 v_h, v_h) + \tau\|\mathbf{s} \cdot \nabla v_h\|^2 - \kappa\tau\|v_h\|\|\mathbf{s} \cdot \nabla v_h\| \\ &\quad - \sigma_s\tau\|S_1 v_h\|\|\mathbf{s} \cdot \nabla v_h\|.\end{aligned}\quad (49)$$

The last two terms can be bounded using Young's inequality and expression (44), for example in the form

$$\begin{aligned}-\kappa\tau\|v_h\|\|\mathbf{s} \cdot \nabla v_h\| - \sigma_s\tau\|S_1 v_h\|\|\mathbf{s} \cdot \nabla v_h\| \\ \geq -\kappa^{1/2}\tau^{1/2}\|v_h\|\|\mathbf{s} \cdot \nabla v_h\| - \sigma_s^{1/2}\tau^{1/2}\|S_1 v_h\|\|\mathbf{s} \cdot \nabla v_h\| \\ \geq -\frac{2}{3}\kappa\|v_h\|^2 - \frac{3}{8}\tau\|\mathbf{s} \cdot \nabla v_h\|^2 - \frac{2}{3}\sigma_s\|S_1 v_h\|^2 - \frac{3}{8}\tau\|\mathbf{s} \cdot \nabla v_h\|^2.\end{aligned}$$

The Lemma follows using the fact that  $-\|S_1 v_h\|^2 \geq -(S_1 v_h, v_h)$ , using the last inequality in (49) and that  $(S_1 v_h, v_h) \geq \lambda_1\|I^{\perp} v_h\|_2$ .  $\square$

**Lemma 3** (Interpolation error of the SUPG method). There holds

$$\mathcal{B}_{\text{supg}}(u - \hat{u}_h, v_h) \lesssim E(h)\|v_h\| \quad \forall v_h \in \mathcal{W}_h.$$

**Proof.** Integrating by parts the first term in  $\mathcal{B}_{\text{supg}}(u - \hat{u}_h, v_h)$  and using Schwartz inequality and the behavior (44) assumed for  $\tau$  it is found that

$$\begin{aligned}\mathcal{B}_{\text{supg}}(u - \hat{u}_h, v_h) &\leq \tau^{-1/2}\|u - \hat{u}_h\|\tau^{1/2}\|\mathbf{s} \cdot \nabla v_h\| + \|u - \hat{u}_h\|_{r^+}\|v_h\|_{r^+} \\ &\quad + \kappa^{1/2}\|u - \hat{u}_h\|\kappa^{1/2}\|v_h\| + \sigma_s^{1/2}\|u - \hat{u}_h\|\sigma_s^{1/2}\|I^{\perp} v_h\| \\ &\quad + \tau^{1/2}\|\nabla(u - \hat{u}_h)\|\tau^{1/2}\|\mathbf{s} \cdot \nabla v_h\| \\ &\quad + \kappa^{1/2}\|u - \hat{u}_h\|\tau^{1/2}\|\mathbf{s} \cdot \nabla v_h\| \\ &\quad + \sigma_s^{1/2}\|u - \hat{u}_h\|\tau^{1/2}\|\mathbf{s} \cdot \nabla v_h\|.\end{aligned}$$

Estimate (48) and the definitions of  $E(h)$  and  $\|\cdot\|$  yield the result.  $\square$

**Theorem 1** (Convergence of the SUPG method). The solution  $u_h$  of problem (23) satisfies

$$\|u - u_h\| \lesssim E(h).$$

**Proof.** The proof is completely standard. From the coercivity given by Lemma 2, the obvious consistency of the SUPG method and the interpolation error estimate in Lemma 3 it follows that

$$\begin{aligned}\|u_h - \hat{u}_h\|^2 &\lesssim \mathcal{B}_{\text{supg}}(u_h - \hat{u}_h, u_h - \hat{u}_h) \\ &= \mathcal{B}_{\text{supg}}(u - \hat{u}_h, u_h - \hat{u}_h) \lesssim E(h)\|u_h - \hat{u}_h\|.\end{aligned}$$

The result is a consequence of this, Lemma 1 and the triangle inequality.  $\square$

#### 4.3. OSS method

In this subsection we prove that method (26) is stable and the solution converges to the continuous one as for the SUPG method. We start proving stability in the form of an inf-sup condition for the bilinear form in (27):

**Lemma 4** (Stability of the OSS method). Suppose that  $c_1$  in (44) is large enough. Then, there is a constant  $C > 0$  such that

$$\inf_{u_h \in \mathcal{W}_h \setminus \{0\}} \sup_{v_h \in \mathcal{W}_h \setminus \{0\}} \frac{\mathcal{B}_{\text{oss}}(u_h, v_h)}{\|u_h\|\|v_h\|} \geq C. \quad (50)$$

**Proof.** Let us start noting that, for any function  $u_h \in \mathcal{W}_h$ , we have

$$\mathcal{B}_{\text{oss}}(u_h, u_h) = \frac{1}{2}\|u_h\|_{r^+}^2 + \kappa\|u_h\|^2 + (S_\sigma u_h, u_h) + \tau\|P_h^\perp(\mathbf{s} \cdot \nabla u_h)\|^2. \quad (51)$$

Clearly,  $\mathcal{B}_{\text{oss}}$  is not coercive in the norm (42). The basic idea is to obtain control on the components on the finite element space for the terms whose orthogonal components appear in this expression. The key point is that this control comes from the Galerkin terms in the bilinear form  $\mathcal{B}_{\text{oss}}$ . Let us consider  $v_{h,0} := \tau P_h(\mathbf{s} \cdot \nabla u_h)$ . We have that

$$\begin{aligned}\mathcal{B}_{\text{oss}}(u_h, v_{h,0}) &= (\mathbf{s} \cdot \nabla u_h + \kappa u_h + S_\sigma u_h, \tau P_h(\mathbf{s} \cdot \nabla u_h)) \\ &\quad + (P_h^\perp(\mathbf{s} \cdot \nabla u_h), \tau \mathbf{s} \cdot \nabla (\tau P_h(\mathbf{s} \cdot \nabla u_h))) \\ &\geq \tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 + \kappa\tau(u_h, P_h(\mathbf{s} \cdot \nabla u_h)) \\ &\quad + \tau(S_\sigma u_h, P_h(\mathbf{s} \cdot \nabla u_h)) \\ &\quad - \tau^2\|P_h^\perp(\mathbf{s} \cdot \nabla u_h)\|\|\nabla(P_h(\mathbf{s} \cdot \nabla u_h))\|.\end{aligned}$$

Using the inverse estimate (45), the fact that  $\tau c_1 \leq h$ , Young's inequality and (11), we get

$$\begin{aligned}\mathcal{B}_{\text{oss}}(u_h, v_{h,0}) &\geq \tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 - \frac{3}{4}\tau\kappa^2\|u_h\|^2 - \frac{3}{4}\tau\sigma_s^2\|S_1 u_h\|^2 \\ &\quad - \frac{2}{3}\tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 - \tau^2 \\ &\quad \times \frac{C_{\text{inv}}}{h}\|P_h^\perp(\mathbf{s} \cdot \nabla u_h)\|\|P_h(\mathbf{s} \cdot \nabla u_h)\| \\ &\geq \frac{1}{3}\tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 - \frac{3}{4}\tau\kappa^2\|u_h\|^2 \\ &\quad - \frac{3}{4}\tau\sigma_s^2\|S_1 u_h\|^2 - \frac{C_{\text{inv}}}{c_1}\tau\|\mathbf{s} \cdot \nabla u_h\|^2.\end{aligned}\quad (52)$$

Let  $v_h = u_h + \alpha v_{h,0}$ . Adding up inequality (52) multiplied by  $\alpha$  to (51) it follows that

$$\begin{aligned}\mathcal{B}_{\text{oss}}(u_h, v_h) &\geq \frac{1}{2}\|u_h\|_{r^+}^2 + \left(1 - \frac{3}{4}\alpha\right)\kappa\|u_h\|^2 + \left(1 - \frac{3}{4}\alpha\right)(S_\sigma u_h, u_h) \\ &\quad + \tau\|P_h^\perp(\mathbf{s} \cdot \nabla u_h)\|^2 + \frac{\alpha}{3}\tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\quad - \frac{C_{\text{inv}}}{c_1}\alpha\tau\|\mathbf{s} \cdot \nabla u_h\|^2.\end{aligned}$$

Using again Young's inequality and that  $\tau\kappa \leq 1$  and  $\tau\sigma_s \leq 1$ , we get

$$\begin{aligned}\mathcal{B}_{\text{oss}}(u_h, v_h) &\geq \frac{1}{2}\|u_h\|_{r^+}^2 + \left(1 - \frac{3}{4}\alpha\right)\kappa\|u_h\|^2 + \left(1 - \frac{3}{4}\alpha\right)(S_\sigma u_h, u_h) \\ &\quad + \tau\|P_h^\perp(\mathbf{s} \cdot \nabla u_h)\|^2 + \frac{\alpha}{3}\tau\|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\quad - \frac{C_{\text{inv}}}{c_1}\alpha\tau\|\mathbf{s} \cdot \nabla u_h\|^2.\end{aligned}$$

Therefore

$$\begin{aligned}\mathcal{B}_{\text{oss}}(u_h, v_h) &\geq \frac{1}{2}\|u_h\|_{r^+}^2 + \left(1 - \frac{3}{4}\alpha\right)\kappa\|u_h\|^2 \\ &\quad + \left(1 - \frac{3}{4}\alpha\right)\lambda_1\sigma_s\|I^{\perp} u_h\|^2 \\ &\quad + \min\left\{1 - \frac{\alpha C_{\text{inv}}}{c_1}, \frac{\alpha}{3} - \frac{\alpha C_{\text{inv}}}{c_1}\right\}\tau\|\mathbf{s} \cdot \nabla u_h\|^2 \\ &\geq \min\left\{1 - \frac{\alpha C_{\text{inv}}}{c_1}, \alpha\left(\frac{1}{3} - \frac{C_{\text{inv}}}{c_1}\right), \left(1 - \frac{3\alpha}{4}\right), \frac{1}{2}\right\}\|u_h\|^2.\end{aligned}\quad (53)$$

If we choose  $\alpha$  such that  $0 < \alpha < \min\left\{\frac{c_1}{C_{\text{inv}}}, \frac{4}{3}\right\}$  we have that  $\mathcal{B}_{\text{oss}}(u_h, v_h) \gtrsim \|u_h\|^2$  for the discrete function  $v_h$  we have chosen, provided  $c_1$  is large enough, for example  $c_1 > 3C_{\text{inv}}$ .

On the other hand, using the inverse and trace inequalities (45) and (47) and condition  $\tau c_1 \leq h$  we have that

$$\begin{aligned} |||v_{h,0}|||^2 &= |||\tau P_h(\mathbf{s} \cdot \nabla u_h)|||^2 \\ &\leq \tau^3 \|\mathbf{s} \cdot \nabla (P_h(\mathbf{s} \cdot \nabla u_h))\|^2 + \tau^2 \|P_h(\mathbf{s} \cdot \nabla u_h)\|_{r^+}^2 + (\kappa \tau^2 \\ &\quad + \lambda_1 \sigma_s \tau^2) \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\leq \tau^3 \frac{C_{\text{inv}}^2}{h^2} \|\mathbf{s} \cdot \nabla u_h\|^2 + \frac{C_{\text{trace}}}{h} \tau^2 \|\mathbf{s} \cdot \nabla u_h\|^2 + 2\tau \|\mathbf{s} \cdot \nabla u_h\|^2 \\ &\leq \left( \frac{C_{\text{inv}}^2}{C_1^2} + \frac{C_{\text{trace}}}{C_1} + 2 \right) \tau \|\mathbf{s} \cdot \nabla u_h\|^2 \lesssim |||u_h|||^2. \end{aligned}$$

Using this fact in (53) we have shown that for each  $u_h \in \mathcal{W}_h$  there exists  $v_h \in \mathcal{W}_h$  such that

$$\mathcal{B}_{\text{oss}}(u_h, v_h) \gtrsim |||u_h||| |||v_h|||,$$

from where the inf sup condition (50) is verified and stability is established.  $\square$

**Lemma 5** (Interpolation error of the OSS method). There holds

$$\mathcal{B}_{\text{oss}}(u - \hat{u}_h, v_h) \lesssim E(h) |||v_h||| \quad \forall v_h \in \mathcal{W}_h.$$

**Proof.** The proof is very similar to that of Lemma 3. The difference is only the treatment of the stabilization term, which in this case can be easily bounded as

$$\begin{aligned} \tau(\mathbf{s} \cdot \nabla(u_h - \hat{u}_h), P_h^\perp(\mathbf{s} \cdot \nabla v_h)) &\leq \tau^{1/2} \|\nabla(u_h - \hat{u}_h)\| \tau^{1/2} \|\mathbf{s} \cdot \nabla v_h\| \\ &\text{from where we can proceed as in Lemma 3. } \square \\ \text{Contrary to the Galerkin and the SUPG methods, the OSS method is not consistent (in the version we have presented it, see [7]). There is a consistency error given by the fact that} \\ \mathcal{B}_{\text{oss}}(u - u_h, v_h) &= \tau(P_h^\perp(\mathbf{s} \cdot \nabla u), \mathbf{s} \cdot \nabla v_h). \end{aligned} \quad (54)$$

However, this consistency error can be bounded as follows:

**Lemma 6** (Bound for the consistency error of the OSS method). Suppose that  $f$  in (1) belongs to the finite element space. Then, there holds

$$\mathcal{B}_{\text{oss}}(u - u_h, v_h) \lesssim E(h) |||v_h||| \quad \forall v_h \in \mathcal{W}_h.$$

**Proof.** From (54) we have that

$$\begin{aligned} \mathcal{B}_{\text{oss}}(u - u_h, v_h) &= \tau(P_h^\perp(\mathbf{s} \cdot \nabla u), \mathbf{s} \cdot \nabla v_h) \\ &\leq \tau^{1/2} \|P_h^\perp(\mathbf{s} \cdot \nabla u)\| \tau^{1/2} \|\mathbf{s} \cdot \nabla v_h\|. \end{aligned}$$

Since  $\mathbf{s} \cdot \nabla u = f - \kappa u - S_\sigma u$  and  $P_h^\perp(f) = 0$ , we have that  $P_h^\perp(\mathbf{s} \cdot \nabla u) = -P_h^\perp(\kappa u + S_\sigma u)$ , and the results follow from the best approximation property of the projection  $P_h$  with respect to the  $L^2$ -norm and the expression of  $\tau$ .  $\square$

**Remark 2.** The assumption  $P_h^\perp(f) = 0$  is not as restrictive as it might seem. Clearly, the component of  $f$  orthogonal to the finite element space vanishes when it is tested with a finite element function, and therefore the Galerkin method does not account for it, in spite of its optimal accuracy (and lack of stability). On the other hand, there would be no problem in keeping the whole residual in the definition of the subscale in Section 3.1.2, case in which the OSS method would be exactly consistent.

Combining the previous results we easily get:

**Theorem 2** (Convergence of the OSS method). The solution  $u_h$  of the OSS method satisfies

$$|||u - u_h||| \lesssim E(h).$$

#### 4.4. Galerkin method

The previous analysis could be easily adapted to account for the possibility  $\tau = 0$ , which corresponds to the Galerkin method. Apart from the need of redefining the error function  $E(h)$  given by (43) if  $\tau = 0$ , also the working norm needs to be modified, since in fact the Galerkin method provides some sort of control on the convective term, as we shall show below. More precisely, we will prove stability and convergence in the norm

$$|||v_h|||_G^2 := \|v_h\|_{r^+}^2 + \kappa \|v_h\|^2 + \lambda_1 \sigma_s \|I\Gamma^\perp v_h\|^2 + h \|P_h(\mathbf{s} \cdot \nabla v_h)\|^2 \quad (55)$$

defined for  $v_h \in \mathcal{W}_h$ . This norm does not contain the whole derivative of  $v_h$  along direction  $\mathbf{s}$ , but only the projection onto the finite element space. Moreover, the factor of the last term is the mesh size  $h$ , not  $\tau$ . We discuss later the implications of these facts.

**Lemma 7** (Stability of the Galerkin method). For  $h(\kappa + \sigma_s) \leq C_0$ , there is a constant  $C > 0$  such that

$$\inf_{u_h \in \mathcal{W}_h \setminus \{0\}} \sup_{v_h \in \mathcal{W}_h \setminus \{0\}} \frac{\mathcal{B}(u_h, v_h)}{|||u_h|||_G |||v_h|||_G} \geq C. \quad (56)$$

**Proof.** For any function  $u_h \in \mathcal{W}_h$ , we have

$$\mathcal{B}(u_h, u_h) = \frac{1}{2} \|u_h\|_{r^+}^2 + \kappa \|u_h\|^2 + (S_\sigma u_h, u_h). \quad (57)$$

It is obvious that  $\mathcal{B}$  is not coercive in the norm (55). Similar to the proof of (50), let us consider  $v_{h,0} := hP_h(\mathbf{s} \cdot \nabla u_h)$ . We have that

$$\begin{aligned} \mathcal{B}(u_h, v_{h,0}) &= (\mathbf{s} \cdot \nabla u_h + \kappa u_h + S_\sigma u_h, hP_h(\mathbf{s} \cdot \nabla u_h)) \\ &= h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 + \kappa h (u_h, P_h(\mathbf{s} \cdot \nabla u_h)) \\ &\quad + h (S_\sigma u_h, P_h(\mathbf{s} \cdot \nabla u_h)). \end{aligned}$$

Using Young's inequality and (11) we can bound this as follows:

$$\begin{aligned} \mathcal{B}(u_h, v_{h,0}) &\geq h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 - \frac{1}{4} h \kappa^2 \|u_h\|^2 \\ &\quad - \frac{3}{4} h \|S_\sigma u_h\|^2 - \frac{2}{3} h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\geq \frac{1}{3} h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 - \frac{3}{4} h \kappa^2 \|u_h\|^2 - \frac{3}{4} h \sigma_s (S_\sigma u_h, u_h). \end{aligned} \quad (58)$$

Thus, from (57) and (58) it follows that, for all  $\alpha > 0$ ,

$$\begin{aligned} \mathcal{B}(u_h, u_h + \alpha v_{h,0}) &\geq \frac{1}{2} \|u_h\|_{r^+}^2 + \left(1 - \frac{3}{4} h \kappa \alpha\right) \kappa \|u_h\|^2 \\ &\quad + \left(1 - \frac{3}{4} h \sigma_s \alpha\right) \lambda_1 \sigma_s \|I\Gamma^\perp u_h\|^2 + \frac{1}{3} \alpha h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2. \end{aligned}$$

From the assumption  $h(\kappa + \sigma_s) \leq C_0$  it follows that we may choose  $\alpha$  such that  $\mathcal{B}(u_h, u_h + \alpha v_{h,0}) \gtrsim |||u_h|||_G^2$ . It remains only to prove that  $|||v_{h,0}|||_G \lesssim |||u_h|||_G$ , which can be done as in the proof of Lemma 4 using (45) and (47) and now condition  $h(\kappa + \sigma_s) \leq C_0$ :

$$\begin{aligned} |||v_{h,0}|||_G^2 &= |||hP_h(\mathbf{s} \cdot \nabla u_h)|||_G^2 \\ &\leq h^3 \|P_h(\mathbf{s} \cdot \nabla (P_h(\mathbf{s} \cdot \nabla u_h)))\|^2 + h^2 \|P_h(\mathbf{s} \cdot \nabla u_h)\|_{r^+}^2 + (\kappa h^2 \\ &\quad + \lambda_1 \sigma_s h^2) \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\leq h^3 \frac{C_{\text{inv}}^2}{h^2} \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 + \frac{C_{\text{trace}}}{h} h^2 \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 + h^2 (\kappa \\ &\quad + \sigma_s) \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \\ &\leq \left( C_{\text{inv}}^2 + C_{\text{trace}} + h(\kappa + \sigma_s) \right) h \|P_h(\mathbf{s} \cdot \nabla u_h)\|^2 \lesssim |||u_h|||_G^2 \end{aligned}$$

from where the result follows.  $\square$



**Table 1**

Convergence behavior of stabilized and Galerkin methods in limiting cases.

Limit case	Stabilized methods	Galerkin method
$\kappa h \ll 1$	$h \ \mathbf{s} \cdot \nabla(u - u_h)\  + (\kappa h)^{1/2} \ u - u_h\  \lesssim h^{p+1}  u _{L^2(S^2; H^{p+1}(\Omega))}$	$h \ P_h(\mathbf{s} \cdot \nabla(u - u_h))\  + (\kappa h)^{1/2} \ u - u_h\  \lesssim (\kappa h)^{-1/2} h^{p+1}  u _{L^2(S^2; H^{p+1}(\Omega))}$
$\kappa h \gg 1$	$h \ \mathbf{s} \cdot \nabla(u - u_h)\  + \kappa h \ u - u_h\  \lesssim (\kappa h) h^{p+1}  u _{L^2(S^2; H^{p+1}(\Omega))}$	$h \ P_h(\mathbf{s} \cdot \nabla(u - u_h))\  + (\kappa h)^{1/2} \ u - u_h\  \lesssim h^{p+1}  u _{L^2(S^2; H^{p+1}(\Omega))}$

**Remark 3.** In stabilized finite element methods one is concerned not only in the asymptotic behavior  $h \rightarrow 0$ , but also in the limits obtained when the physical parameters vary for  $h$  fixed. Assumption  $h(\kappa + \sigma_s) \leq C_0$  precludes estimate (56) to be valid when  $h$  is fixed and either  $\kappa \rightarrow \infty$  or  $\sigma_s \rightarrow \infty$ . However, this is a theoretical restriction rather than a practical one in the problem we are analyzing. In particular, it does not appear in the numerical examples of Section 5. In fact, when  $\kappa$  is constant, as we assume, it can be relaxed to  $h\sigma_s \leq C_0$  for stability (Lemma 7) but condition  $h(\kappa + \sigma_s) \leq C_0$  will be again required for convergence (Lemma 8 and Theorem 3). To see why  $h\sigma_s \leq C_0$  is enough for stability it suffices to observe that

$$\kappa h(u_h, P_h(\mathbf{s} \cdot \nabla u_h)) = \kappa h(u_h, \mathbf{s} \cdot \nabla u_h) = \kappa h \frac{1}{2} \|u_h\|_{L^2}^2 \geq 0,$$

so that bound (58) is not sharp. However, the proof of Lemma 7 can be straightforwardly extended to variable coefficients if condition  $h(\kappa + \sigma_s) \leq C_0$ , with  $\kappa$  the maximum value of the absorption coefficient, is kept.

The error function of the Galerkin method is determined by the following result:

**Lemma 8** (Interpolation error of the Galerkin method). For  $h(\kappa + \sigma_s) \leq C_0$ , there holds

$$\|u - \tilde{u}_h\|_G \lesssim h^{1/2} \varepsilon_1(u), \quad (59)$$

$$B(u - u_h, v_h) \lesssim (\kappa^{-1/2} + h^{1/2}) \varepsilon_1(u) \|v_h\|_G. \quad (60)$$

**Proof.** Estimate (59) is a trivial consequence of the definition of  $\|\cdot\|_G$  in (55), assumption  $h(\kappa + \sigma_s) \leq C_0$ , the trace inequality (47) and the fact that  $h^{1/2} \varepsilon_1(u) = h^{-1/2} \varepsilon_0(u)$ . The proof of (60) is as follows:

$$\begin{aligned} B(u - \tilde{u}_h, v_h) &= (\mathbf{s} \cdot \nabla(u - \tilde{u}_h), v_h) + \kappa(u - \tilde{u}_h, v_h) + (S_\sigma(u - \tilde{u}_h), v_h) \\ &\leq \|\nabla(u - \tilde{u}_h)\| \|v_h\| + \kappa \|u - \tilde{u}_h\| \|v_h\| \\ &\quad + \sigma_s \|u - \tilde{u}_h\| \|\Pi^\perp v_h\| \lesssim (\kappa^{-1/2} \varepsilon_1(u) \\ &\quad + \kappa^{1/2} \varepsilon_0(u)) \kappa^{1/2} \|v_h\| + \sigma_s^{1/2} \varepsilon_0(u) \sigma_s^{1/2} \|\Pi^\perp v_h\| \\ &\lesssim [\kappa^{-1/2} + (\kappa^{1/2} + \sigma_s^{1/2})h] \varepsilon_1(u) \|v_h\|_G. \end{aligned} \quad (61)$$

The proof is complete using once again that  $h(\kappa + \sigma_s) \leq C_0$ .  $\square$

Combining the results of Lemmas 7 and 8 and the consistency of the Galerkin formulation it is found that:

**Theorem 3** (Convergence of the Galerkin method). For  $h(\kappa + \sigma_s) \leq C_0$ , the solution  $u_h$  of the Galerkin method satisfies

$$\|u - u_h\|_G \lesssim (\kappa^{-1/2} + h^{1/2}) \varepsilon_1(u).$$

At this point it is interesting to compare what happens in the limit of dominant directional derivative or dominant absorption in the stabilized formulations, either SUPG or OSS, and the Galerkin method. To simplify the discussion, suppose that  $\sigma_s = 0$  and let us neglect the error control obtained on the boundary  $\Gamma^+$ . We may write the convergence estimates as

Stabilized methods(SUPG, OSS) :

$$\begin{aligned} \tau^{1/2} \|\mathbf{s} \cdot \nabla(u - u_h)\| + \kappa^{1/2} \|u - u_h\| \\ \lesssim (\tau^{-1/2} h^{p+1} + \tau^{1/2} h^p) |u|_{L^2(S^2; H^{p+1}(\Omega))}, \end{aligned}$$

Galerkin method :

$$\begin{aligned} h^{1/2} \|P_h(\mathbf{s} \cdot \nabla(u - u_h))\| + \kappa^{1/2} \|u - u_h\| \\ \lesssim (\kappa^{-1/2} h^p + h^{1/2} h^p) |u|_{L^2(S^2; H^{p+1}(\Omega))}. \end{aligned}$$

The behavior when  $\kappa h \ll 1$  (small absorption) and when  $\kappa h \gg 1$  (large absorption) is displayed in Table 1 (results in this table are obtained multiplying by adequate factors both sides of the corresponding error estimates). The two main conclusions that may be drawn from this table are:

- When absorption is dominant, both stabilized methods and the Galerkin method yield optimal convergence in the  $L^2(\Omega)$ -norm of the error. Note however that we cannot consider  $h$  fixed and let  $\kappa \rightarrow \infty$  because of the assumption  $h(\kappa + \sigma_s) \leq C_0$  on which all our previous analysis relies (see also Remark 3). Thus, the estimate for  $\kappa h \gg 1$  in the case of the Galerkin method has to be understood with caution, considering that  $\kappa h$  is large but without the possibility to take the limit  $\kappa h \rightarrow \infty$ . It cannot be considered better than the estimate for the stabilized methods.
- When absorption is small, stabilized methods yield optimal convergence, of order  $h^p$  for the directional derivative. However, the Galerkin method fails because of the large factor  $(\kappa h)^{-1/2}$  in the estimate. This is due to the fact that the norm of the directional derivative is controlled).

## 5. Numerical examples

To investigate and compare the accuracy and efficiency of the SUPG and OSS stabilization methods, three typical test problems with absorbing/emitting and scattering media enclosed by gray walls are considered. The directional domain  $S^2$  is discretized with the discrete ordinates method, using the  $S_N$  quadrature sets introduced by Lathrop and Carlson [14]. Three tests cases are selected to compare the behavior of the different methods.

After the spatial and directional discretizations have been carried out, the resulting linear system of equations is expensive and strongly coupled due to the discretization of the integral operator in the RTE (1). In order to save computer memory, the DOM equations are solved iteratively. If  $l$  denotes the iteration counter, the implemented iterative scheme is

$$\begin{aligned} \mathbf{s} \cdot \nabla u_l^\alpha + \left( \kappa + \sigma_s - \frac{\sigma_s}{4\pi} w_\alpha \phi(\mathbf{s}^\alpha, \mathbf{s}^\alpha) \right) u_l^\alpha \\ = \frac{\sigma_s}{4\pi} \sum_{\beta=1, \beta \neq \alpha}^N w_\beta \phi(\mathbf{s}^\alpha, \mathbf{s}^\beta) u_{l-1}^\beta + f^\alpha, \end{aligned} \quad (62)$$

with  $\alpha = 1, 2, \dots, N$ , and where  $\mathbf{s}^\alpha$  and  $w_\alpha$  are the chosen sets of directions (ordinates) and weights; the unknown  $u_l^\alpha$  is the radiative intensity propagating in direction  $\mathbf{s}^\alpha$  evaluated at iteration  $l$ , and  $f^\alpha$  is the source term. We have to deal with  $N$  equations that are solved independently for each direction, and that are coupled only by the

right hand side. The scheme described is of Jacobi type, although a Gauss–Seidel iterative method could also be employed.

In some of the cases described next, the radiative transfer equation (1) is subject to emissive and reflective boundary conditions of the form (12). As explained in Section 2.3, the variational problem (13) needs to be discretized. In the stabilized formulation, the forms involved in the problem are modified as explained for homogeneous boundary conditions. After the DOM discretization, the new boundary conditions couple different directions through the reflective integral term. In the framework of the iterative scheme (62), values of  $u$  at the previous iteration can be used to evaluate the resulting right-hand-side term.

Let us consider the finite element approximation of (62). For each equation of this system the implementation is based on an *a priori* calculation of the integrals appearing in the formulation and then the construction of the matrices and right-hand-side vectors of the final algebraic systems to be solved. These matrices and these vectors can be constructed directly for each nodal point, without the need to loop over the elements, thus making the calculations much faster.

It is important to note that as  $(\kappa + \sigma_s)h \rightarrow 0$  each discrete equation is dominated by the convective term. The Galerkin contribution of the convective term is a singular matrix. Therefore, as  $(\kappa + \sigma_s)h \rightarrow 0$  the Galerkin method gives almost singular matrices. This causes that iterative solvers as GMRES do not converge, even when using good ILUT preconditioners. We had to use direct solvers for solving the test problems described next with the Galerkin method.

### 5.1. Gaussian shaped radiative source term between one-dimensional parallel black slabs

The first test problem that we consider, taken from [21], is known to produce Galerkin oscillations. It consists in solving the radiative transfer problem in a nonscattering medium between

one-dimensional finite parallel black slab. This problem is modeled by the one dimensional RTE

$$\mu \frac{du}{dx} + \kappa u = \exp(-2500(x - 0.5)^2), \quad x \in [0, 1],$$

where  $\mu$  is the cosine between direction  $\mathbf{s}$  and the  $x$  axis.

Homogeneous boundary conditions are taken:

$$u(0, \mu) = 0, \quad \mu > 0,$$

$$u(1, \mu) = 0, \quad \mu < 0.$$

The analytical solution of this problem in the case of  $\mu > 0$  can be written as

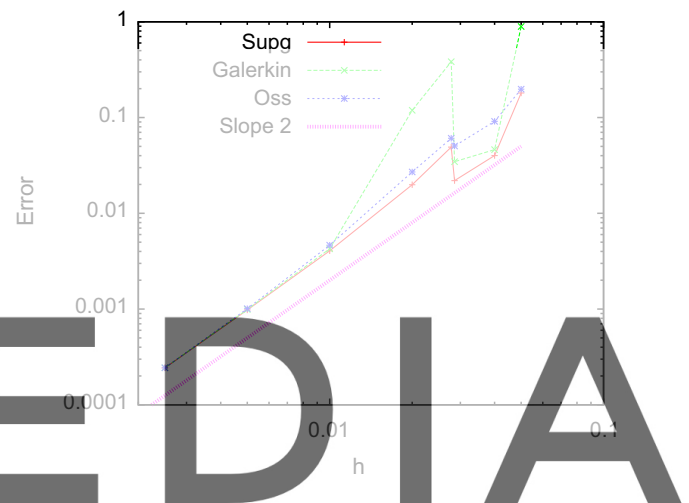


Fig. 2. Relative error of the solutions for different methods against mesh size  $h$ .

Register for free at <https://www.scipedia.com> to download the version without the watermark

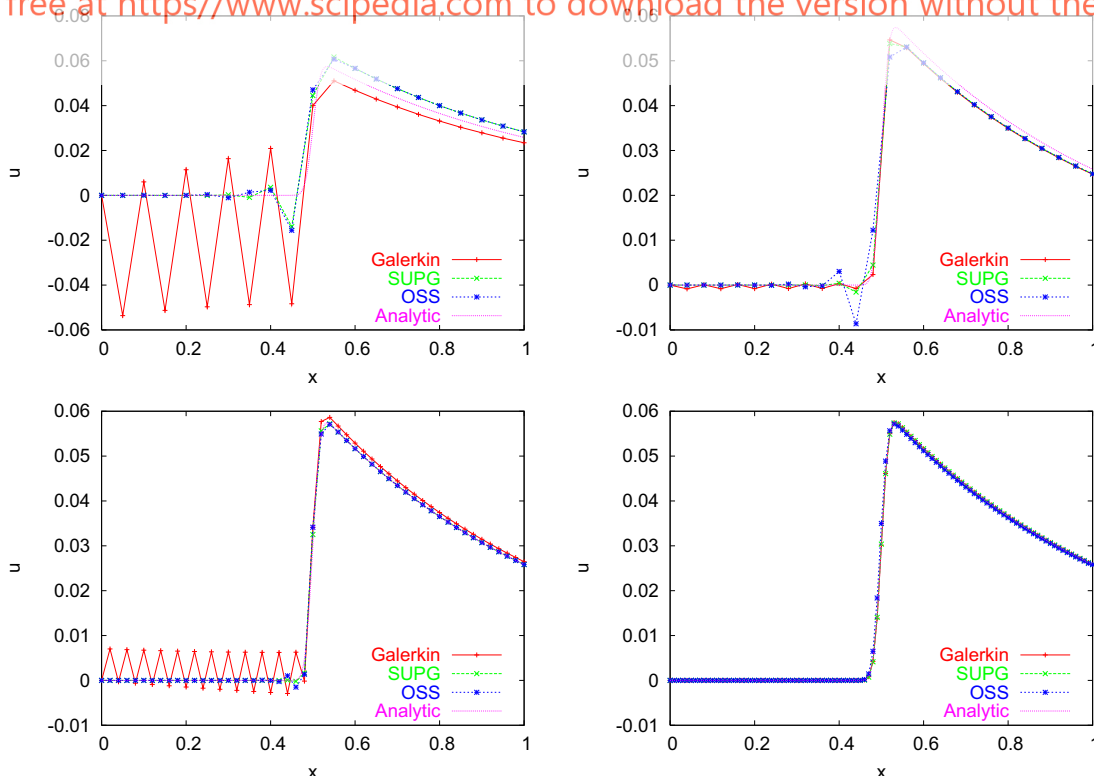


Fig. 1. Radiative intensity distribution for a mesh with 20 (top left), 25 (top right), 50 (bottom left) and 100 (bottom right) elements.

$$u(x) = -\frac{0.02\sqrt{\pi}}{2\mu} \exp\left(-\frac{\kappa}{\mu}\left(x - \frac{\kappa}{10000\mu} - 0.5\right)\right) \times \left(\operatorname{erf}\left(\frac{\kappa}{100\mu} + 50(0.5 - x)\right) - \operatorname{erf}\left(\frac{\kappa}{100\mu} + 25\right)\right).$$

Fig. 1 shows the radiative intensity distribution for  $\mu = 0.5773505$  and  $\kappa = 1$ . The Galerkin and the stabilized SUPG and OSS methods are compared against the analytical solution. We used uniform meshes, with the number of elements ranging from 20 to 400 linear elements. For coarser grids global spurious oscillations occur when the Galerkin method is used. This is because it is stable in the norm (55), which has poor control on the derivatives. Due to the nature of the analytical solution we found bigger oscillations for an even quantity of elements.

When using stabilized formulations global oscillations are removed. For some meshes, the OSS method presents higher localized peaks than the SUPG method due to the less diffusive nature of the scheme. When using a finer grid capable of capturing the jump of the analytical solution, all oscillations are removed.

Fig. 2 shows the  $L^2$  error of the different methods relative to the reference solution against the mesh size  $h$ . We observe from this figure optimal convergence, that is  $\|u - u_h\| \leq Ch^2$  when  $h \rightarrow 0$ .

## 5.2. Absorbing and anisotropic scattering in the unit square (2D problem)

The second test problem consists in solving the radiative heat transfer equation over a square domain. The medium is considered

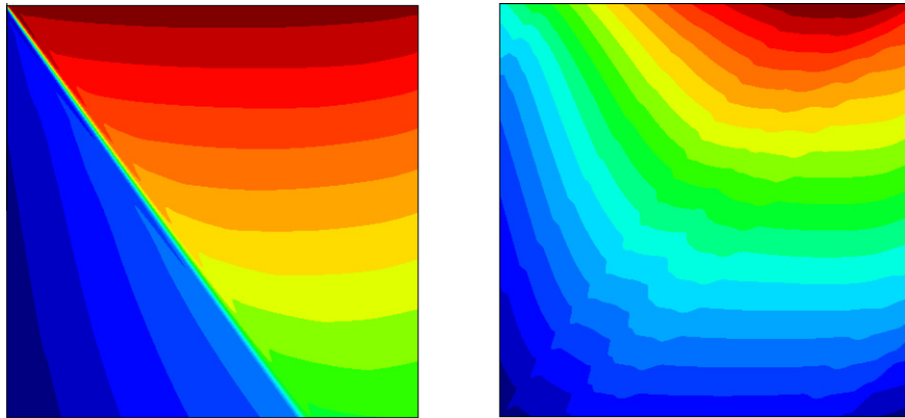


Fig. 3. Radiative intensity solutions using the SUPG method in a mesh of  $240 \times 240$  elements corresponding to case 1 ( $\kappa = 0.2 \text{ m}^{-1}$  and  $\sigma_s = 0.8 \text{ m}^{-1}$ ). Radiation propagating from the upper hot wall (left picture) and the cold walls (right picture).

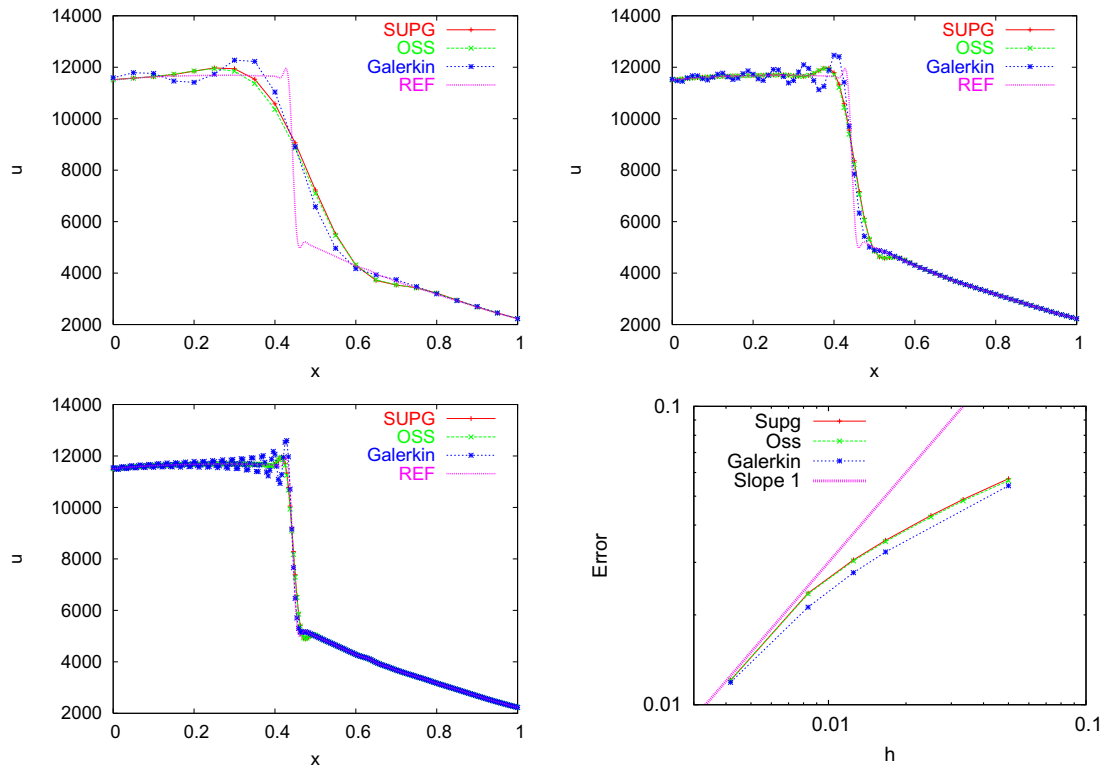


Fig. 4. Radiative intensity cut at  $x = 0.4 \text{ m}$  using different numerical methods on a mesh of  $20 \times 20$  (top left),  $80 \times 80$  (top right) and  $240 \times 240$  (bottom left) elements. In all cases the reference solution computed on a mesh of  $480 \times 480$  elements is also shown. Relative error of the solutions against mesh size  $h$  (bottom right). Case 1 ( $\kappa = 0.2$ ).

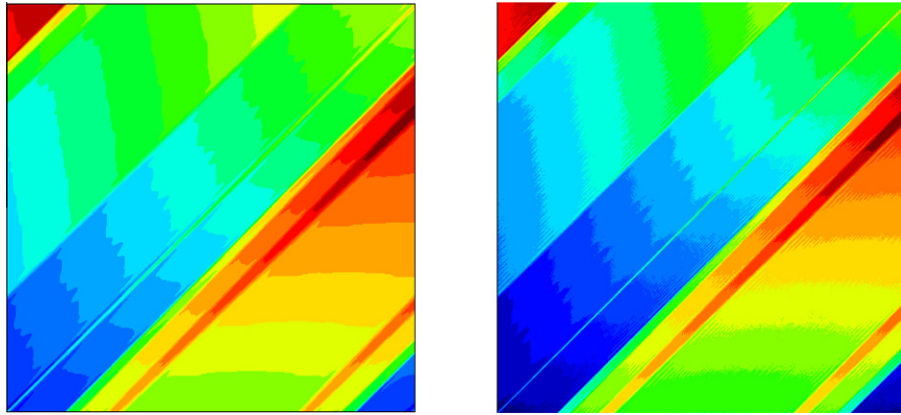


Fig. 5. Radiative intensity solution using the SUPG (left) and the Galerkin (right) methods in a mesh of  $480 \times 480$  elements. Case 2 ( $\kappa = 0.01 \text{ m}^{-1}$  and  $\sigma_s = 0.001 \text{ m}^{-1}$ ).

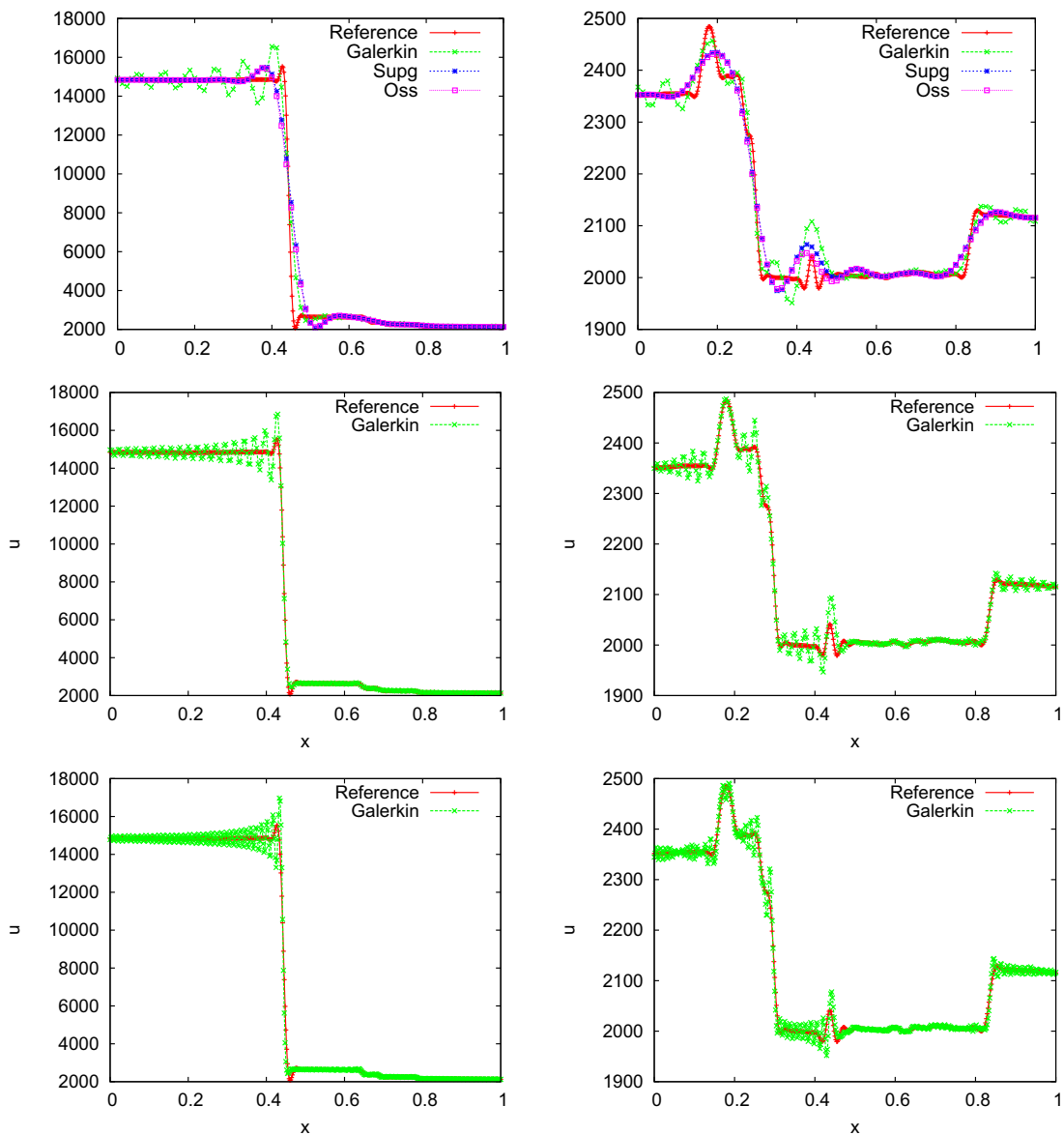


Fig. 6. Radiative intensity cut at 0.4 m (right) and 0.6 m (left) from the hot wall for different numerical methods and reference solution in meshes of  $80 \times 80$  (top)  $240 \times 240$  (middle)  $480 \times 480$  (bottom) elements. Case 2 ( $\kappa = 0.01 \text{ m}^{-1}$  and  $\sigma_s = 0.001 \text{ m}^{-1}$ ).



to be absorbing/emitting and with anisotropic scattering, enclosed by boundaries of length  $L = 1$  m with emissivity  $\epsilon = 0.8$  and reflectivity  $\rho = 0.2$ . From the discussion after (15) it follows that the maximum upper bound for  $\eta$  in this case is  $\eta \approx 49.49$ . We have chosen  $\eta = 40$  with good results.

The upper wall is maintained at a temperature  $T_{\text{hot}} = 1000$  K, and all other walls at a temperature  $T_{\text{cold}} = 500$  K. The medium is maintained at uniform temperature of  $T_g = 800$  K. We consider the phase function  $\phi(\mathbf{s}', \mathbf{s})$  as linearly anisotropic, of the form

$$\phi(\mathbf{s}', \mathbf{s}) = 1 + A_1 \mathbf{s}' \cdot \mathbf{s},$$

where  $A_1 = 0.2$ .

For the space discretization we use bilinear rectangular elements  $Q_1$ . For discretizing the angular dependency is used the SN8 ordinates set [16], consisting of 80 directions. Discretization in space goes from meshes of  $20 \times 20$  to  $240 \times 240$  elements.

As the RTE does not have analytical solutions for arbitrary geometries, we have compared the results of the different formulations with a reference solution, obtained using a finer grid of  $480 \times 480$  elements.

We have run two cases, the first one with an optical thickness  $\beta = (\kappa + \sigma_s)L = 1$ , and the second one with  $\beta = (\kappa + \sigma_s)L = 0.011$ . In the later case the medium is quite nonparticipative. Due to the discontinuity in the boundary conditions, solutions may present sharp gradients that can activate instabilities.

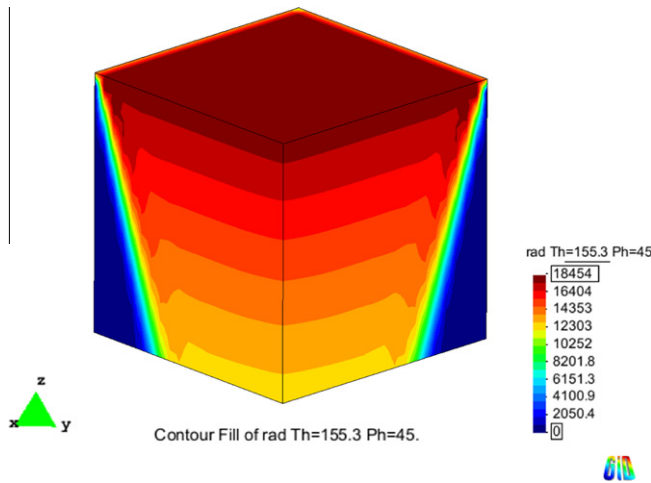
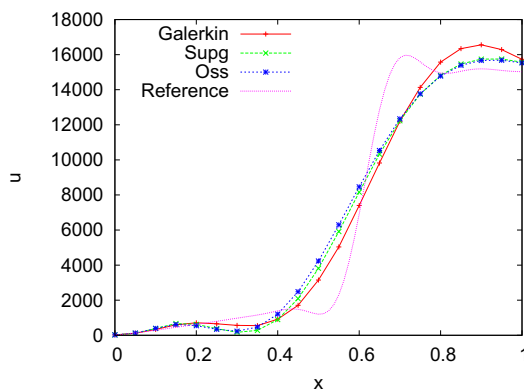


Fig. 7. Radiative intensity solution with the SUPG method in the whole domain using a mesh of  $20 \times 20 \times 20$  trilinear elements.



**Case 1** In this first case the physical properties of the medium are taken as  $\kappa = 0.2 \text{ m}^{-1}$  and  $\sigma_s = 0.8 \text{ m}^{-1}$ . In Fig. 3, two solutions are plotted for the radiation intensity in  $\Omega$  when solving with the SUPG method and using a mesh of  $240 \times 240$  elements. The picture on the left corresponds to radiation propagating from the upper hot wall, whereas in the picture on the right it propagates from the cold walls.

The solutions obtained when using Galerkin, the SUPG and the OSS methods are compared against the reference solution in Fig. 4, using meshes of  $20 \times 20$ ,  $80 \times 80$  and  $240 \times 240$  linear elements. A cut of these solutions is shown. The Galerkin method shows higher numerical oscillations for finer grids. The stabilized methods OSS and SUPG give very similar results, without numerical oscillations.

The error of the different methods relative to reference solution is also plotted against mesh size  $h$  in Fig. 4 (bottom-right). This error has been computed as

$$\text{Error} = \frac{\sum_{a,\alpha} (u_h(\mathbf{x}^a, \mathbf{s}^\alpha) - u(\mathbf{x}^a, \mathbf{s}^\alpha))^2}{\sum_{a,\alpha} (u(\mathbf{x}^a, \mathbf{s}^\alpha))^2}, \quad (63)$$

where  $a, \alpha$  refers to nodes and directions,  $u_h(\mathbf{x}^a, \mathbf{s}^\alpha)$  is the discrete solution at node  $\mathbf{x}^a$  and direction  $\mathbf{s}^\alpha$ , and  $u(\mathbf{x}^a, \mathbf{s}^\alpha)$  is the reference solution at this node and with this direction. For the smallest values of  $h$  a linear convergence of the error is observed. This convergence is not optimal. A possible explanation is that the analytical solution is discontinuous due to the discontinuity in the boundary conditions.

**Case 2** In this case the medium has as absorption coefficient  $\kappa = 0.01 \text{ m}^{-1}$  and as scattering coefficient  $\sigma_s = 0.001 \text{ m}^{-1}$ . Fig. 5 shows the solutions obtained for the radiation intensity when using the SUPG and the Galerkin methods in a mesh of  $480 \times 480$  elements. In this example, the radiation intensity comes from the upper hot wall. It is observed that the Galerkin solution is polluted with global oscillations.

In Fig. 6 different cuts of radiative intensity are shown for the stabilized and the Galerkin methods for meshes of  $80 \times 80$ ,  $240 \times 240$  and  $480 \times 480$  elements. As in case 1, the Galerkin method shows higher numerical oscillations for finer grids. When using finer grids, the OSS and the SUPG methods give results similar to the reference solution, so that only this reference solution has been plotted. It is worth to note that the Galerkin oscillations are not node to node.

### 5.3. Absorbing and isotropic scattering in the unit cube (3D problem)

The third test problem consists in solving the radiative transfer equation in the unit cube  $(x, y, z) \in [0, 1]^3$ . The temperature of the medium is  $T_m = 800$  K. The boundary conditions consist of one hot wall ( $z = 1$ ) at  $T_h = 1000$  K, while the other walls are maintained

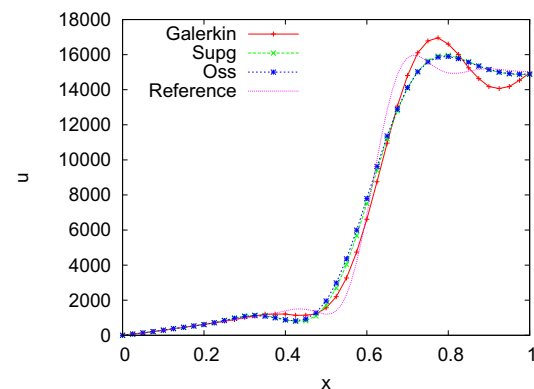


Fig. 8. Radiative intensity cut for the SUPG, the OSS and the Galerkin methods against the reference solution in a mesh of  $20 \times 20 \times 20$  (left) and  $40 \times 40 \times 40$  (right) trilinear elements.

cold ( $T_c = 0$  K). The hot wall is considered opaque and non reflective ( $\epsilon = 1.0, \rho = 0.0$ ). The physical absorption and scattering are  $\kappa = 0.2$  and  $\sigma_s = 0.3$  (SI units are assumed).

Fig. 7 shows the solution for radiation intensity over the cubic domain coming from the upper hot wall when using the SUPG method and a mesh of  $20 \times 20 \times 20$  trilinear elements. Fig. 8 shows plots of radiative intensity cuts for the stabilized and the Galerkin methods using respectively meshes of  $20 \times 20 \times 20$  and  $40 \times 40 \times 40$  trilinear elements. The tests are compared to a reference solution obtained with a mesh of  $80 \times 80 \times 80$  trilinear elements using the SUPG method. The obtained results using the SUPG and the OSS methods are very close. The behavior of the Galerkin method is similar to the one observed in the bidimensional problem (see Fig. 4).

## 6. Conclusions

In this paper we have designed and analyzed stabilized finite element methods to approximate the radiative transport equation. The problem is posed in a spatial domain and in the unit sphere  $S^2$ , and both need to be discretized. We have focused our attention to the spatial discretization, and used only the DOM in the numerical testing, although any other discretization of  $S^2$  could be used.

The Galerkin method for the spatial discretization suffers from numerical oscillations due to the convective term in the equation to be approximated. We have presented a non-conventional numerical analysis that shows that some control on the convective derivative can be obtained, but not enough to prevent the appearance of numerical wiggles.

In order to overcome the misbehavior of the Galerkin method, two stabilized finite element methods have been discussed, namely, the well known SUPG formulation and the OSS method. Both can be motivated within the variational multiscale framework, although some simplifying assumptions have to be added to arrive to the version of the methods analyzed here.

Both approximations, the SUPG and the OSS, are stable and optimally convergent in the same norm and with the same error function. This norm happens to be finer than the one in which the Galerkin method can be analyzed. There is full control in the convective derivative that translates into globally smooth solutions, although some local oscillations may be still encountered. As the OSS method introduces less numerical dissipation than the SUPG method, the local overshoots and undershoots are sometimes higher using the OSS method. Let us stress that the norm in which we have presented the stability and convergence results remains meaningful for all values of the physical parameters. For the SUPG method this represents a modification of well known results, whereas for the OSS method this analysis was not available.

## Acknowledgment

This work has been partially supported by project FITUN, Ref. TRA2008-05162, from the Spanish Ministry of Science and Innovation. Financial support provided by the International Center For Mechanical Sciences (CISM) and Autostrada del Brennero S.p.A., from Italy, is also acknowledged.

## References

- [1] M. Asadzadeh, A finite element method for the neutron transport equation in an infinite cylindrical domain, *SIAM J. Numer. Anal.* 35 (1998) 1299–1314.
- [2] F. Brezzi, D. Marini, E. Süli, Residual-free bubbles for advection–diffusion problems: the general error analysis, *Numer. Math.* 85 (2000) 31–47.
- [3] A.N. Brooks, T.J.R. Hughes, Streamline upwind/Petrov–Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier–Stokes equations, *Comput. Meth. Appl. Mech. Engrg.* 32 (1982) 199–259.
- [4] S. Chandrasekhar, *Radiative Transfer*, Dover Publications, 1960.
- [5] R. Codina, Stabilization of incompressibility and convection through orthogonal sub-scales in finite element methods, *Comput. Meth. Appl. Mech. Engrg.* 190 (13–14) (2000) 1579–1599.
- [6] R. Codina, Stabilized finite element approximation of transient incompressible flows using orthogonal subscales, *Comput. Meth. Appl. Mech. Engrg.* 191 (2002) 4295–4321.
- [7] R. Codina, Analysis of a stabilized finite element approximation of the oseen equations using orthogonal subscales, *Appl. Numer. Math.* 58 (2008) 264–283.
- [8] R. Codina, Finite element approximation of the hyperbolic wave equation in mixed form, *Comput. Meth. Appl. Mech. Engrg.* 197 (2008) 1305–1322.
- [9] R. Codina, Finite element approximation of the three field formulation of the Stokes problem using arbitrary interpolations, *SIAM J. Numer. Anal.* 47 (2009) 699–718.
- [10] A. Ern, J.L. Guermond, Discontinuous Galerkin methods for Friedrichs's systems: I. general theory, *SIAM J. Numer. Anal.* 44 (2006) 753–778.
- [11] T.J.R. Hughes, G.R. Feijóo, L. Mazzei, J.B. Quincy, The variational multiscale method – a paradigm for computational mechanics, *Comput. Meth. Appl. Mech. Engrg.* 166 (1998) 3–24.
- [12] J.H. Jeans, The equations of radiative transfer of energy, *Mon. Not. R. Astron. Soc.* 78 (1917) 28–36.
- [13] G. Kanschat, Robust finite element discretization for radiative transfer problems with scattering, *East-West J. Numer. Math.* 6 (4) (1998) 265–272.
- [14] K.D. Lathrop, B.G. Carlson, Discrete-ordinates angular quadrature of the neutron transport equation, Technical Information Series Report LASL-3186, Los Alamos Scientific Laboratory, 1965.
- [15] M.M. Razzaque, D.E. Klein, J.R. Howell, Finite element solution of radiative heat transfer in two dimensional rectangular enclosure with gray participating media, *J. Heat Transf.* 105 (1983) 933–936.
- [16] M.F. Modest, *Radiative Heat Transfer*, Academic Press, 2003.
- [17] J.A. Nitsche, Über ein variationsprinzip zur lösung von dirichlet-problemen bei verwendung von teilräumen, die keinen randbedingungen unterworfen sind, *Abhandlungen aus dem Mathematischen Seminar der Universität Hamburg* 36 (1971) 9–15.
- [18] R. Viskanta, Heat transfer by conduction and radiation in absorbing and scattering materials, *J. Heat Transf.* 143 (1965) 143–150.
- [19] S. Richling, E. Meinköhn, N. Kryzhevoi, G. Kanschat, Radiative transfer with finite elements, *Astron. Astrophys.* 380 (2001) 776–788.
- [20] W.A. Fiveland, Finite element formulation of the discrete ordinates method for multidimensional geometries, *J. Thermophys. Heat Transf.* 8 (1994) 426–433.
- [21] J.M. Zhao, L.H. Liu, Second order radiative transfer equation and its properties of numerical solution using the finite element method, *Numer. Heat Transf. Part B* 51 (2007) 391–409.